

09/701626

Practitioner's Docket No. NEB-165-PUS

## CHAPTER II

## Preliminary Classification:

Proposed Class:

Subclass:

NOTE: "All applicants are requested to include a preliminary classification on newly filed patent applications. The preliminary classification, preferably class and subclass designations, should be identified in the upper right-hand corner of the letter of transmittal accompanying the application papers, for example 'Proposed Class 2, subclass 129.'" M.P.E.P., § 601, 7th ed.

**TRANSMITTAL LETTER  
TO THE UNITED STATES ELECTED OFFICE (EO/US)**

## (ENTRY INTO U.S. NATIONAL PHASE UNDER CHAPTER II)

PCT/US99/13295

11 June 1999

12 June 1998

INTERNATIONAL APPLICATION NO.

INTERNATIONAL FILING DATE

PRIORITY DATE CLAIMED

Restriction Enzyme Gene Discovery Method

TITLE OF INVENTION

Elisabeth A. Raleigh, Romualdas Vaisvila, Richard D. Morgan

APPLICANT(S)

## Box PCT

Assistant Commissioner for Patents

Washington D.C. 20231

ATTENTION: EO/US

## CERTIFICATION UNDER 37 C.F.R. § 1.10\*

(Express Mail label number is mandatory.)

(Express Mail certification is optional.)

I hereby certify that this Transmittal Letter and the papers indicated as being transmitted therewith is being deposited with the United States Postal Service on this date 1 December 2000, in an envelope as "Express Mail Post Office to Addressee" Mailing Label Number EL010489946US, addressed to the Assistant Commissioner for Patents, Washington, D.C. 20231.

Melissa A. Jackson

(type or print name of person mailing paper)

Signature of person mailing paper

**WARNING:** Certificate of mailing (first class) or facsimile transmission procedures of 37 C.F.R. § 1.8 cannot be used to obtain a date of mailing or transmission for this correspondence.

**\*WARNING:** Each paper or fee filed by "Express Mail" **must** have the number of the "Express Mail" mailing label placed thereon prior to mailing. 37 C.F.R. § 1.10(b).

"Since the filing of correspondence under § 1.10 without the Express Mail mailing label thereon is an oversight that can be avoided by the exercise of reasonable care, requests for waiver of this requirement will **not** be granted on petition." Notice of Oct. 24, 1996, 60 Fed. Reg. 56,439, at 56,442.

(Transmittal Letter to the United States Elected Office (EO/US) [13-18]—page 1 of 8)

001001 09/701626

NOTE: To avoid abandonment of the application, the applicant shall furnish to the USPTO, not later than 30 months from the priority date: (1) a copy of the international application, if it has been previously communicated by the International Bureau or unless it was originally filed in the USPTO; and (2) the basic national fee (see 37 C.F.R. § 1.492(a)). The 30-month time limit may not be extended. 37 C.F.R. § 1.495.

**WARNING:** Where the items are those which can be submitted to complete the entry of the international application into the national phase are subsequent to 30 months from the priority date the application is still considered to be in the international state and if mailing procedures are utilized to obtain a date the express mail procedure of 37 C.F.R. § 1.10 must be used (since international application papers are not covered by an ordinary certificate of mailing—See 37 C.F.R. § 1.8.

NOTE: Documents and fees must be clearly identified as a submission to enter the national state under 35 U.S.C. § 371 otherwise the submission will be considered as being made under 35 U.S.C. § 111. 37 C.F.R. § 1.494(f).

- I. Applicant herewith submits to the United States Elected Office (EO/US) the following items under 35 U.S.C. § 371:
- ☒ This express request to immediately begin national examination procedures (35 U.S.C. § 371(f)).
  - ☒ The U.S. National Fee (35 U.S.C. § 371(c)(1)) and other fees (37 C.F.R. § 1.492) as indicated below:

(Transmittal Letter to the United States Elected Office (EO/US) [13-18]—page 2 of 8)

## 2. Fees

CLAIMS FEE	(1) FOR	(2) NUMBER FILED	(3) NUMBER EXTRA	(4) RATE	(5) CALCULATIONS
<input checked="" type="checkbox"/> *	TOTAL CLAIMS	17 - 20 =	0	× \$18.00 =	\$ 0.00
	INDEPENDENT CLAIMS	5 - 3 =	2	× \$78.00 =	156.00
	MULTIPLE DEPENDENT CLAIM(S) (if applicable) + \$260.00				0.00
	<b>BASIC FEE**</b> <input checked="" type="checkbox"/> U.S. PTO WAS INTERNATIONAL PRELIMINARY EXAMINATION AUTHORITY Where an International preliminary examination fee as set forth in § 1.482 has been paid on the international application to the U.S. PTO: <input type="checkbox"/> and the international preliminary examination report states that the criteria of novelty, inventive step (non-obviousness) and industrial activity, as defined in PCT Article 33(1) to (4) have been satisfied for all the claims presented in the application entering the national stage (37 C.F.R. § 1.492(a)(4)) ..... \$96.00 <input checked="" type="checkbox"/> and the above requirements are not met (37 C.F.R. § 1.492(a)(1)) ..... \$670.00 <input type="checkbox"/> U.S. PTO WAS NOT INTERNATIONAL PRELIMINARY EXAMINATION AUTHORITY Where no international preliminary examination fee as set forth in § 1.482 has been paid to the U.S. PTO, and payment of an international search fee as set forth in § 1.445(a)(2) to the U.S. PTO: <input type="checkbox"/> has been paid (37 C.F.R. § 1.492(a)(2)) ..... \$690.00 <input type="checkbox"/> has not been paid (37 C.F.R. § 1.492(a)(3)) ..... \$970.00 <input type="checkbox"/> where a search report on the international application has been prepared by the European Patent Office or the Japanese Patent Office (37 C.F.R. § 1.492(a)(5) ) ..... \$840.00				670.00
Total of above Calculations					= 826.00
SMALL ENTITY	Reduction by 1/2 for filing by small entity, if applicable. Affidavit must be filed also. (note 37 C.F.R. § 1.9, 1.27, 1.28)				- 428.00
	Subtotal				428.00
	Total National Fee				\$ 428.00
	Fee for recording the enclosed assignment document \$40.00 (37 C.F.R. § 1.21(h)). (See Item 13 below). See attached "ASSIGNMENT COVER SHEET".				40.00
TOTAL	Total Fees enclosed				\$ 468.00

\*See attached Preliminary Amendment Reducing the Number of Claims.

09/701626

i. ☒ A check in the amount of 468.00 to cover the above fees is enclosed.

ii. ☐ Please charge Account No. \_\_\_\_\_ in the amount of \$ \_\_\_\_\_.

A duplicate copy of this sheet is enclosed.

**\*\*WARNING:** "To avoid abandonment of the application the applicant shall furnish to the United States Patent and Trademark Office not later than the expiration of 30 months from the priority date: \* \* \* (2) the basic national fee (see § 1.492(a)). The 30-month time limit may not be extended." 37 C.F.R. § 1.495(b).

**WARNING:** If the translation of the international application and/or the oath or declaration have not been submitted by the applicant within thirty (30) months from the priority date, such requirements may be met within a time period set by the Office. 37 C.F.R. § 1.495(b)(2). The payment of the surcharge set forth in § 1.492(e) is required as a condition for accepting the oath or declaration later than thirty (30) months after the priority date. The payment of the processing fee set forth in § 1.492(f) is required for acceptance of an English translation later than thirty (30) months after the priority date. Failure to comply with these requirements will result in abandonment of the application. The provisions of § 1.136 apply to the period which is set. Notice of Jan. 3, 1993, 1147 O.G. 29 to 40.

3. ☒ A copy of the International application as filed (35 U.S.C. § 371(c)(2)):

**NOTE:** Section 1.495 (b) was amended to require that the basic national fee and a copy of the international application must be filed with the Office by 30 months from the priority date to avoid abandonment. "The International Bureau normally provides the copy of the international application to the Office in accordance with PCT Article 20. At the same time, the International Bureau notifies applicant of the communication to the Office. In accordance with PCT Rule 47.1, that notice shall be accepted by all designated offices as conclusive evidence that the communication has duly taken place. Thus, if the applicant desires to enter the national stage, the applicant normally need only check to be sure the notice from the International Bureau has been received and then pay the basic national fee by 30 months from the priority date." Notice of Jan. 7, 1993, 1147 O.G. 29 to 40, at 35-36. See item 14c below.

a. ☐ is transmitted herewith.

b. ☒ is not required, as the application was filed with the United States Receiving Office.

c. ☐ has been transmitted

i. ☐ by the International Bureau.

Date of mailing of the application (from form PCT/1B/308): \_\_\_\_\_

ii. ☐ by applicant on \_\_\_\_\_

Date

4. ☒ A translation of the International application into the English language (35 U.S.C. § 371(c)(2)):

a. ☐ is transmitted herewith.

b. ☒ is not required as the application was filed in English.

c. ☐ was previously transmitted by applicant on \_\_\_\_\_

Date

d. ☐ will follow.

(Transmittal Letter to the United States Elected Office (EO/US) [13-18]—page 4 of 8)

5. ☒ Amendments to the claims of the International application under PCT Article 19 (35 U.S.C. § 371(c)(3)):

NOTE: The Notice of January 7, 1993 points out that 37 C.F.R. § 1.495(a) was amended to clarify the existing and continuing practice that PCT Article 19 amendments must be submitted by 30 months from the priority date and this deadline may not be extended. The Notice further advises that: "The failure to do so will not result in loss of the subject matter of the PCT Article 19 amendments. Applicant may submit that subject matter in a preliminary amendment filed under section 1.121. In many cases, filing an amendment under section 1.121 is preferable since grammatical or idiomatic errors may be corrected." 1147 O.G. 29-40, at 36.

- a. ☐ are transmitted herewith.
- b. ☐ have been transmitted
  - i. ☐ by the International Bureau.  
Date of mailing of the amendment (from form PCT/1B/308): \_\_\_\_\_
  - ii. ☐ by applicant on (date) \_\_\_\_\_  
Date
- c. ☒ have not been transmitted as
  - i. ☒ applicant chose not to make amendments under PCT Article 19.  
Date of mailing of Search Report (from form PCT/ISA/210): 19 Oct. 1999
  - ii. ☐ the time limit for the submission of amendments has not yet expired.  
The amendments or a statement that amendments have not been made will be transmitted before the expiration of the time limit under PCT Rule 46.1.

6. ☒ A translation of the amendments to the claims under PCT Article 19 (38 U.S.C. § 371(c)(3)):

- a. ☐ is transmitted herewith.
- b. ☒ is not required as the amendments were made in the English language.
- c. ☐ has not been transmitted for reasons indicated at point 5(c) above.

7. ☒ A copy of the international examination report (PCT/IPEA/409)

- ☐ is transmitted herewith.
- ☒ is not required as the application was filed with the United States Receiving Office.

8. ☒ Annex(es) to the international preliminary examination report

- a. ☐ is/are transmitted herewith.
- b. ☒ is/are not required as the application was filed with the United States Receiving Office.

9. ☒ A translation of the annexes to the international preliminary examination report

- a. ☐ is transmitted herewith.
- b. ☒ is not required as the annexes are in the English language.

10. ☒ An oath or declaration of the inventor (35 U.S.C. § 371(c)(4)) complying with 35 U.S.C. § 115

a. ☐ was previously submitted by applicant on \_\_\_\_\_

Date

b. ☒ is submitted herewith, and such oath or declaration

i. ☒ is attached to the application this transmittal

ii. ☐ identifies the application and any amendments under PCT Article 19 that were transmitted as stated in points 3(b) or 3(c) and 5(b); and states that they were reviewed by the inventor as required by 37 C.F.R. § 1.70.

c. ☐ will follow.

II. Other document(s) or information included:

11. ☐ An International Search Report (PCT/ISA/210) or Declaration under PCT Article 17(2)(a):

a. ☐ is transmitted herewith.

b. ☐ has been transmitted by the International Bureau.  
Date of mailing (from form PCT/IB/308): \_\_\_\_\_

c. ☐ is not required, as the application was searched by the United States International Searching Authority.

d. ☐ will be transmitted promptly upon request.

e. ☐ has been submitted by applicant on \_\_\_\_\_

Date

12. ☐ An Information Disclosure Statement under 37 C.F.R. §§ 1.97 and 1.98:

a. ☐ is transmitted herewith.

Also transmitted herewith is/are:

☐ Form PTO-1449 (PTO/SB/08A and 08B).

☐ Copies of citations listed.

b. ☐ will be transmitted within THREE MONTHS of the date of submission of requirements under 35 U.S.C. § 371(c).

c. ☐ was previously submitted by applicant on \_\_\_\_\_

Date

13. ☒ An assignment document is transmitted herewith for recording.

A separate ☐ "COVER SHEET FOR ASSIGNMENT (DOCUMENT) ACCOMPANYING NEW PATENT APPLICATION" or ☒ FORM PTO 1595 is also attached.

New England Biolabs, Inc.

32 Tozer Road

Beverly, MA 01915

(Transmittal Letter to the United States Elected Office (EO/US) [13-18]—page 6 of 8)

09/701626

14. ☒ Additional documents:

- a. ☐ Copy of request (PCT/RO/101)
- b. ☐ International Publication No. \_\_\_\_\_
- i. ☐ Specification, claims and drawing
- ii. ☐ Front page only
- c. ☒ Preliminary amendment (37 C.F.R. § 1.121)
- d. ☒ Other

Sequence Listing in computer-readable format, paper copy  
and statement regarding the same

15. ☒ The above checked items are being transmitted

- a. ☒ before 30 months from any claimed priority date.
- b. ☐ after 30 months.

16. ☐ Certain requirements under 35 U.S.C. § 371 were previously submitted by the applicant on \_\_\_\_\_, namely:

---



---



---



---



---

**AUTHORIZATION TO CHARGE ADDITIONAL FEES**

**WARNING:** Accurately count claims, especially multiple dependant claims, to avoid unexpected high charges if extra claims are authorized.

**NOTE:** "A written request may be submitted in an application that is an authorization to treat any concurrent or future reply, requiring a petition for an extension of time under this paragraph for its timely submission, as incorporating a petition for extension of time for the appropriate length of time. An authorization to charge all required fees, fees under § 1.17, or all required extension of time fees will be treated as a constructive petition for an extension of time in any concurrent or future reply requiring a petition for an extension of time under this paragraph for its timely submission. Submission of the fee set forth in § 1.17(a) will also be treated as a constructive petition for an extension of time in any concurrent reply requiring a petition for an extension of time under this paragraph for its timely submission." 37 C.F.R. § 1.136(a)(3).

**NOTE:** "Amounts of twenty-five dollars or less will not be returned unless specifically requested within a reasonable time, nor will the payer be notified of such amounts; amounts over twenty-five dollars may be returned by check or, if requested, by credit to a deposit account." 37 C.F.R. § 1.26(a).

- ☒ The Commissioner is hereby authorized to charge the following additional fees that may be required by this paper and during the entire pendency of this application to Account No. 14-0740.

- ☒ 37 C.F.R. § 1.492(a)(1), (2), (3), and (4) (filing fees)

**WARNING:** Because failure to pay the national fee within 30 months without extension (37 C.F.R. § 1.495(b)(2)) results in abandonment of the application, it would be best to always check the above box.

☐ 37 C.F.R. § 1.492(b), (c) and (d) (presentation of extra claims)

NOTE: Because additional fees for excess or multiple dependent claims not paid on filing, or on later presentation, must only be paid or these claims cancelled by amendment prior to the expiration of the time period set for response by the PTO in any notice of fee deficiency (37 C.F.R. § 1.492(d)), it might be best not to authorize the PTO to charge additional claim fees, except possible when dealing with amendments after final action.

☐ 37 C.F.R. § 1.17 (application processing fees)

☐ 37 C.F.R. § 1.17(a)(1)-(5) (extension fees pursuant to § 1.136(a).

☐ 37 C.F.R. § 1.18 (issue fee at or before mailing of Notice of Allowance, pursuant to 37 C.F.R. § 1.311(b))

NOTE: Where an authorization to charge the issue fee to a deposit account has been filed before the mailing of a Notice of Allowance, the issue fee will be automatically charged to the deposit account at the time of mailing the notice of allowance. 37 C.F.R. § 1.311(b).

NOTE: 37 C.F.R. § 1.28(b) requires "Notification of any change in loss of entitlement to small entity status must be filed in the application . . . prior to paying, or at the time of paying . . . issue fee." From the wording of 37 C.F.R. § 1.28(b): (a) notification of change of status must be made even if the fee is paid as "other than a small entity" and (b) no notification is required if the change is to another small entity.

☐ 37 C.F.R. § 1.492(e) and (f) (surcharge fees for filing the declaration and/or filing an English translation of an International Application later than 30 months after the priority date).

Reg. No.: 30901

Tel. No.: (978 ) 927-5054 X:292

Customer No.:

  
SIGNATURE OF PRACTITIONER

Gregory D. Williams

General Counsel

(type or print name of practitioner)

New England Biolabs, Inc.

P.O. Address

32 Tozer Road

Beverly, MA 01915

(Transmittal Letter to the United States Elected Office (EO/US) [13-18]—page 8 of 8)



Docket: NEB-165-PUS

IN THE UNITED STATES ELECTED OFFICE (EO/US)

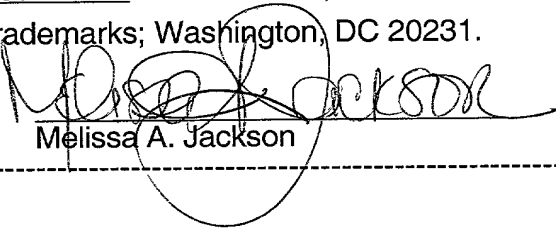
International Application No.: PCT/US99/13295  
International Filing Date: 11 June 1999  
Priority Date Claimed: 12 June 1998  
Title of Invention: Restriction Enzyme Gene Discovery Method  
Applicant(s): Raleigh, et al..  
Box PCT  
Commissioner of Patents  
and Trademarks  
Washington, DC 20231

---

I, Melissa A. Jackson hereby certify that the following documents are being deposited, via Express Mail, on this date, 1st ~~November~~ <sup>11</sup> December 2000:

1. Transmittal Letter to the United States Elected Office (Entry Into U.S. National Phase Under Chapter II);
2. Recordation of Assignment;
3. Assignment;
4. Declaration and Power of Attorney;
5. Sequence Listing (disk), Papercopy;
6. Statement regarding Submission;
7. Preliminary Amendment; and
8. Check in the amount of \$468.00

in an envelope addressed as "Express Mail Post Office to Addressee"  
Mailing Label Number EL010489946US to: BOX PCT; Honorable  
Commissioner of Patents and Trademarks; Washington, DC 20231.

  
Melissa A. Jackson

---

Sir:

**PRELIMINARY AMENDMENT**

Applicants wish to amend the above-identified Published Application as follows:

Parameter	Value	Unit
Initial concentration	0.1	mol/L
Temperature	25	°C
Time	0, 1, 2, 4, 8, 16, 32, 64, 128, 256, 512, 1024, 2048, 4096, 8192, 16384, 32768, 65536, 131072, 262144, 524288, 1048576, 2097152, 4194304, 8388608, 16777216, 33554432, 67108864, 134217728, 268435456, 536870912, 1073741824, 2147483648, 4294967296, 8589934592, 17179869184, 34359738368, 68719476736, 137438953472, 274877906944, 549755813888, 1099511627776, 2199023255552, 4398046511104, 8796093022208, 17592186044416, 35184372088832, 70368744177664, 140737488355328, 281474976710656, 562949953421312, 1125899906842624, 2251799813685248, 4503599627370496, 9007199254740992, 18014398509481984, 36028797018963968, 72057594037927936, 144115188075855872, 288230376151711744, 576460752303423488, 1152921504606846976, 2305843009213693952, 4611686018427387904, 9223372036854775808, 18446744073709551616, 36893488147419103232, 73786976294838206464, 147573952589676412928, 295147905179352825856, 590295810358705651712, 1180591620717411303424, 2361183241434822606848, 4722366482869645213696, 9444732965739290427392, 18889465931478580854784, 37778931862957161709568, 75557863725914323419136, 151115727451828646838272, 302231454903657293676544, 604462909807314587353088, 1208925819614629174706176, 2417851639229258349412352, 4835703278458516698824704, 9671406556917033397649408, 19342813113834066795298816, 38685626227668133590597632, 77371252455336267181195264, 154742504910672534362390528, 309485009821345068724781056, 618970019642690137449562112, 1237940039285380274899124224, 2475880078570760549798248448, 4951760157141521099596496896, 9903520314283042199192993792, 19807040628566084398385987584, 39614081257132168796771975168, 79228162514264337593543950336, 158456325028528675187087900672, 316912650057057350374175801344, 633825300114114700748351602688, 1267650600228229401496703205376, 2535301200456458802993406410752, 5070602400912917605986812821504, 10141204801825835211973625643008, 20282409603651670423947251286016, 40564819207303340847894502572032, 81129638414606681695789005144064, 162259276829213363391578010288128, 324518553658426726783156020576256, 649037107316853453566312041152512, 1298074214633706907132624082305024, 2596148429267413814265248164610048, 5192296858534827628530496329220096, 10384593717069655257060992658440192, 20769187434139310514121985316880384, 41538374868278621028243970633760768, 83076749736557242056487941267521536, 166153499473114484112975882535043072, 332306998946228968225951765070086144, 664613997892457936451903530140172288, 1329227995784915872903807060280344576, 2658455991569831745807614120560689152, 5316911983139663491615228241121378304, 10633823966279326983230456482242756608, 21267647932558653966460912964485513216, 42535295865117307932921825928971026432, 85070591730234615865843651857942052864, 170141183460469231731687303715884105728, 340282366920938463463374607431768211456, 680564733841876926926749214863536422912, 1361129467683753853853498429727072845824, 2722258935367507707706996859454145691648, 5444517870735015415413993718908291383296, 10889035741470030830827987437816582766592, 21778071482940061661655974875633165533184, 43556142965880123323311949751266331066368, 87112285931760246646623899502532662132736, 174224571863520493293247799005065324265472, 348449143727040986586495598010130648530944, 696898287454081973172991196020261297061888, 1393796574908163946345982392040522594123776, 2787593149816327892691964784081045188247552, 5575186299632655785383929568162090376495104, 11150372599265311570767859136324180752990208, 22300745198530623141535718272648361505980416, 44601490397061246283071436545296723011960832, 89202980794122492566142873090593446023921664, 178405961588244985132285746181186892047843328, 356811923176489970264571492362373784095686656, 713623846352979940529142984724747568191373312, 1427247692705959881058285969449495136382746624, 2854495385411919762116571938898990272765493248, 5708990770823839524233143877797980545530986496, 11417981541647679048466287755595961091061972992, 228359630832953580969325755111919221821239	

At page 46, line 5, replace "on \_\_\_\_\_ 1999 and has received ATCC Patent Deposit No. \_\_\_\_\_" with --on June 11, 1999 and has received ATCC Deposit No. PTA-215--.

**REMARKS**

Applicants have amended the specification, specifically page 46, lines 5 and 10-11 to incorporate the ATCC Deposit information which was unavailable at the time of the Application was filed. No new matter has been added by virtue of the amendments made to the specification.

It is respectfully requested that these amendments be entered in the above-identified PCT Application.

Respectfully submitted,

NEW ENGLAND BIOLABS, INC.

Gregory D. Williams  
(Reg. No. 30901)  
New England Biolabs, Inc.  
32 Tozer Road  
Beverly, Massachusetts 01915  
(508) 927-5054; Ext. 292

## RESTRICTION ENZYME GENE DISCOVERY METHOD

### RELATED APPLICATIONS

5 This Application is a PCT Application of U.S. Provisional Application Serial No. 60/089,101 filed 12 June 1998 and U.S. Provisional Application Serial No. 60/089,086 filed 12 June 1998, the disclosures of which are hereby incorporated by reference herein.

### FIELD OF THE INVENTION

10 The invention is generally directed to the field of gene discovery, cloning and expression. A particular aspect of the invention is that it enables direct cloning of intact genes, with a high probability that the orientation of expression is known in advance, and with a low probability of being associated with extraneous possibly toxic genes

15 The invention is limited to genes of a particular kind, since some genes are more likely to be susceptible to cloning and discovery by this method than other genes. Accordingly, the invention is more specifically directed to cloning of genes found within arrays of gene cassettes separated by conserved repeated sequences. Based on present understanding, such arrays are found in prokaryotic organisms and contain genes that have functions that are selectively advantageous to a high level under certain circumstances but are not required under other conditions.

20 Accordingly, some kinds of genes will not be found within these arrays, while other kinds of genes should be enriched in such arrays. Among the genes to be found in such cassette arrays are many genes of commercial interest. The kinds of genes of interest that may be expected in such arrays include:

25  
30 Restriction enzymes, which are useful for a variety of procedures in molecular biology and which enable construction of many useful vectors.

Adhesins, which may allow a cell to attach to a particular surface. Enabling specific attachment to a particular surface rather than others has many uses in providing coatings and targeting molecules or organisms to locations of interest. Such adhesins may also mediate pathogenic processes when expressed by pathogenic organisms, and availability of an adhesin may enable competitive exclusion of such pathogenic organisms.

Small-molecule modifying enzymes, which may convert a toxic or other material abundant in a particular environment to another less toxic to humans or animals, or into a form more useful.

Specific toxin molecules that interact with a host organism, which may be useful for synthesis of inhibitors or antagonists of the toxin or for vaccine purposes.

Different examples of related cassette-encoded gene products will have common general properties (adhesins stick to things) but highly variable specificities (there are many different kinds of specific surfaces to stick to, from rocks to intestinal mucosa to urinary epithelium). Genes of this kind will be referred to below as "diversity-selected genes". The list of gene types above is not exhaustive.

## BACKGROUND OF THE INVENTION

### Hypervariable gene regions in prokaryotic organisms

Hypervariable regions, which show a high level of sequence divergence between closely related strains of the same species, are found at various positions in prokaryotic chromosomes. In some cases, genes present in one strain are absent entirely from a close relative. Examples of this phenomenon include so-called "pathogenicity islands", chromosomal elements that carry genes required for pathogenesis (McDaniel, et al., *Proc. Natl Acad. Sci. USA* 92(5):1664-1668 (1995)). Restriction enzyme genes are sometimes found in regions that are hypervariable in this way (Daniel, et al., *J. Bacteriol.* 170:1775-1782 (1988); Raleigh, *Mol. Microbiol.* 6:1079-1086 (1992); Barcus, et al., *Genetics*, 140:1187-

1197 (1995)). The mechanism of assembly and variation of these regions may depend on novel genetic mechanisms.

### **Integrans and superintegrans as hypervariable gene regions: mobile gene cassettes**

Integrans (Hall and Collis, *Mol. Microbiol.*, 15(4):593-600 (1995)) are arrays of promoterless gene cassettes, separated by related DNA elements ("59 bp elements") that are sites of action for site-specific integrases related to the lambda integrase (Fig. 1). Each integran has at the 5' end a gene for the relevant integrase. Within the integrase gene is a promoter oriented toward the cassettes, upon which expression of all cassette-borne genes is dependent. Cassettes can be found as extrachromosomal nonreplicating circles, and these can be inserted into the array by the integrase. Characterized integrans are plasmid-borne, and the cassettes specify resistance to drugs or other toxic products (such as mercury). Ordinary integrans are small: up to 8 cassettes have been identified in one ordinary integran, and most have between one and three. It is thought that all the genes are expressed from the single promoter found within the sequence of the flanking integrase (Levesque, et al., *Gene* 142(1):49-54 (1994); Recchia and Hall, *Mol. Microbiol.*, 15(1):179-187 (1995)) (Fig. 1); in any event, promoter-like sequences are usually not identified within the gene cassettes. The plasmid location and the multiple-drug resistant character of integrans probably reflect the historical origins of the studies involved: they were found as a result of studies on horizontal transmission of drug resistance in bacteria isolated from clinical settings, where such behavior is selectively advantageous.

A superintegrin (Mazel, et al., *Science*, 280(5363):605-608 (1998)) was recently described as a chromosomal array of a large number of gene cassettes mobilizable by a site-specific integrase obtained from an integrin. This large array, found in *Vibrio cholerae*, may contain up to a hundred cassettes and may account for as much as 10% of the chromosome (Barker, et al., *J. Bacteriol.*, 176(17) 5450-5458 (1994)). The Manning laboratory identified this array in the course of studying a pathogenesis-related hemagglutinin (Franzon, et al., *Infect. Immun.*, 61(7):3032-3037 (1993)). Open reading frames within this array are separated by repeated sequences called VCR (for *Vibrio cholerae* repeats). These repeats are similar to but not the same as the "59 bp elements" of drug-resistance

integrans (Mazel, et al., *supra* (1998)). Manning's laboratory claims to have identified an integrase associated with *Vibrio cholerae* (Clark, et al., *Mol. Microbiol.*, 26(5):1137-1138 (1997)), and the Davies laboratory has published a description of such a gene from *Vibrio cholerae* (Mazel, et al., *supra* (1998)).

This superintegron is distinguished from the ordinary integrons in four respects: size, placement of promoters, replicon location, and the nature of the genes found within cassettes. In contrast to the best-studied integron examples, there appear to be 60 to 100 cassettes within the *V. cholerae* array; and since they are not all oriented in the same direction (Fig. 2), they cannot be expressed from a common promoter. Moreover, the functions encoded by the superintegron are apparently diverse, and some are possibly related to pathogenesis (Mazel, et al., *supra* (1998)). Some of the cassette-borne genes were related to some plasmid-encoded proteins (from database-matching of ORFs 3.1 and 3.2 of the sequence reported in (Barker, et al., *supra* (1994)), one was a heat-stable toxin (Ogawa and Takeda, *Microbiol. Immunol.* 37(8):607-616 (1993)), and one was similar to a lipoprotein gene (*vlpA*: from database matching of ORF2). Accordingly, we surmise (following Mazel et al) that this array may function to cluster genes related to pathogenicity and to the entrap genes specifying other biochemical functions.

### Repeated sequences between gene cassettes in integrons and superintegrons

The sequences interspersed between gene cassettes are thought to be responsible for acquisition and exchange of gene cassettes among the various replicons on which they are located. These sequences, designated "59 bp elements" or "VCR elements" are diverse in sequence but display some common features. A consensus sequence was initially deduced for conventional "59 bp elements" (Hall, et al., *Mol. Microbiol.*, 5(8):1941-1959 (1991)), consisting of:

5' GYCTAACAA-TTCGTTCAAGCCGACGCCGC-T...

ICS

...-TC-GCGGC-GCGGCTTAAGTC-ARGCGTTAGRY 3' (SEQ ID NO:92)

CS

It was later found that the relevant sequences varied in length and sequence within the segments (Hall and Collis, *supra* (1995)). Two most conserved segments could always be identified: 5' to a gene cassette (and at the 3' end of the sequence above; underlined) is found the "Core Sequence" (CS), GTT<sup>rrrr</sup>RY (SEQ ID NO:93); and 3' to a cassette (and at the 5' end of the sequence above; underlined) is found the "Inverse Core Sequence" (ICS), RY<sup>yy</sup>YAAC (SEQ ID NO:94). These two elements are related as inverted repeats. Upon excision, the part of the sequence included in the extrachromosomal circle includes the sequence 3' to the gene as far as the G in the Core Sequence; the circle is completed with the remainder of the CS from the 5' end of the gene (TTAGRY (SEQ ID NO:95)).

The VCR elements were originally said to be unrelated to any other sequence (Barker, et al., *supra* (1994)) but were subsequently shown to conform with the specifications of the "59 bp elements" except for greater length (Mazel, et al., *supra* (1998); Clark, et al., *supra* (1997)): they consist of 124-bp direct repeats of imperfect dyad symmetry, and carry ICS and CS motifs at the ends. VCR elements were found nine times in the original sequence surrounding the putative hemagglutinin gene (Barker, et al., *supra* (1994)).

PCR has been used for characterization of integrons. Some studies employed primers annealing to the conserved integrase genes, or to *sull*, a conserved gene found at the 3' end of many integrons (e.g. (Levesque, et al., *Antimicrob. Agents Chemother.*, 39(1):185-191 (1995); Sallen, et al., *Microb. Drug Resist.*, 1(3):195-202 (1995); Sandvang, et al., *FEMS Microbiol. Lett.*, 160(1):37-41 (1998)). Other studies have employed primers annealing to particular cassette-encoded genes (e.g. (Senda, et al., *J. Clin. Microbiol.* 35(12):2909-2913 (1996); Tosini, et al., *Antimicrob. Agents Chemother.*, 42(12):3053-3058 (1998)). However, it has been considered unlikely that these repeat sequences would enable acquisition of cassette-encoded genes by PCR, because of the degeneracy of the sequences and the secondary structure encoded by them (Hall and Stokes, *Genetica*, 90(2-3):115-132 (1993)). Mazel et al (*supra*, (1998)) were able to obtain cassettes by PCR using primers annealing to the VCR elements, however.

## Background of restriction enzyme gene discovery

Restriction enzyme properties.

Restriction enzymes are the workhorses of molecular biology research. They specifically recognize sites in DNA of 4 to 8 basepairs in length, with extremely high selectivity--that is, a site with one mismatch is typically recognized with an affinity one-thousandfold less than the affinity shown for the correct site. This high degree of selectivity is essential for use in practical applications.

Known restriction enzymes recognize over 200 different specific DNA sequences (Roberts and Macelis, *Nucleic Acids Res.*, 26(1):338-350 1998)) and many of these are commercially available. However, the potential number of different sites is much larger: 32,512 distinct 8-base sites might be recognized  $(((4^8/2)-256)$ : a site 8 bases in length with four possible bases at each position; which can be recognized in either of two complementary strands; minus 256, since 8-base palindromes each read the same in the two strands].

Enzymes with 8 bp recognition sites (8-cutters, such as NotI, SfiI, SmaI, PacI and PmeI) are of particular utility. These enzymes are used for constructing maps of and manipulating DNA from high-complexity sources, such as the genomes of humans and other higher eukaryotes. This utility arises from the rarity of the sites (once per 65,000 bp for palindromic sites), enabling for example the isolation of a whole gene with large introns on a single DNA fragment.

Of the twelve known specificities with 8 bp recognition sites, two were found in *Pseudomonas spp.*, nine in *Streptomyces* or other high G+C gram positive bacteria, and one in *Staphylococcus*. Sequence information is available for six of these, the two *Pseudomonas* isolates and four from high G+C organisms.

Competing approaches to restriction enzyme discovery.

In the past, two broad approaches have been taken to the problem of finding new restriction enzymes: screening for new enzymatic activities, and changing existing enzymes to recognize new sites.



1) Screening of crude extracts of individual prokaryotic strains (obtained from strain collections or natural environments). A test substrate (e.g. phage lambda DNA) is incubated with such an extract, and the digest visualized by agarose gel electrophoresis. This standard approach identifies at least one site-specific nuclease in about 25% of crude extracts screened, with the routine use of targets of combined complexity of about 200 kb.

This approach has two critical defects. First, the fraction of such enzymes recognizing new sites is now very low. In part this may be due to its bias toward identifying enzymes with recognition sites between four and six bp in length and inefficiency in detecting enzymes with larger targets, which are frequently not present in the target substrates.

The second defect is that is extremely labor-intensive. Each strain must be examined individually, and several of the steps involved are projects in themselves: culture growth, cell lysis, and extract clarification each can be a custom procedure. The quality of crude extract preparations varies greatly among isolates, in the extent of contamination with extraneous nucleases, DNA binding proteins and proteases.

In the specific case of *Pseudomonas* and its relatives, extracts are frequently difficult to handle due to extensive nuclease contamination. *Xanthomonas* strains (which are relatives of *Pseudomonas*) frequently give cultures that are hard to collect by centrifugation due to copious extracellular polysaccharide production, and extracts are difficult to clarify for the same reason.

2) Mutational alteration of existing enzymes so that they recognize new sequences. Starting with enzymes recognizing 6 base pairs for which structural information is available, attempts have been made to alter specificity by site-directed, random or random cassette mutagenesis (e.g. (Dorner and Schildkraut, *Nucleic Acids Res.* 22(6):1068-1074 (1994); Heitman and Model, *EMBO J.* 9(10):3369-3378 (1990); Ivanenko, et al., *Biol. Chem.* 379(4-5):459-465 (1998); Hager, et al., *J. Biol. Chem.* 265(35):21520-21526 (1990) and I. Schildkraut, personal communication). Although this work may eventually yield useful products, it has not yet produced an increased specificity (recognizing more bases) or altered specificity (recognizing a different sequence of the same length).

## Background of restriction enzyme gene clone identification and cloning

5        Restriction enzymes are found in a wide variety of prokaryotic organisms, many of them with fastidious growth requirements and frequently in low amounts. For purposes of commercial production, it is most useful to be able to produce a restriction enzyme in a well-understood and genetically tractable bacterial host such as *Escherichia coli*. The many tools for gene expression and regulation, as well as for genetic manipulation of the host cell, enable preparations to be made with  
10       higher purity and lower cost. Accordingly it is very useful to obtain the genes for endonucleases as molecular clones.

### Methyltransferase selection method

15       One method for identifying the presence of a restriction enzyme gene in a clone library is to rely on the presence and expression of a closely-linked gene for a cognate DNA methyltransferase (Wilson, U.S. Patent No. 5,200,333 (1993)). Such methyltransferase enzymes recognize specific DNA sequences and add a methyl group to an A or C residue within the sequence. This modification prevents cleavage by the endonuclease, thereby protecting the host genome from lethal  
20       damage. If such a methyltransferase gene is present in a clone library and effectively expressed, the DNA of that clone will be protected from digestion. This enables selection for the clone in vitro: plasmid clone DNA is purified from a pool of clones and digested with the desired endonuclease enzyme. The methyltransferase clone will not be digested, while other clones in the library,  
25       (which are found in different cells) will be destroyed. Retransformation following such a procedure allows establishment of a selected pool, in which representation of the methyltransferase gene is greatly enriched. If the endonuclease gene is adjacent to the methyltransferase gene, as is often the case, then that gene (or a portion of it) will also be recovered frequently. This method is called the  
30       "methyltransferase selection" method. It is quite useful when three conditions obtain: a cognate methyltransferase exists; the genes for the two functions are tightly linked in the DNA; and the methyltransferase is expressed in *E. coli*.

Several modifications have been added to this basic method, enabling isolation of the endonuclease gene when the first clone does not contain the complete endonuclease gene or when the methyltransferase must be expressed in the cell first, before the endonuclease can be introduced (the "two-step" method) (Brooks and Howard, U.S. Patent No.5,320,957 (1994)).

#### Degenerate methyltransferase-motif PCR method

A second method for identifying the presence of a restriction system gene pair in a clone library is to rely on the presence of conserved polypeptide motif elements found in the DNA methyltransferase proteins (Klimasauskas, et al., *Nucleic Acids Res.* 17:9823-9832 (1989); Lauster, et al., *J. Mol. Biol.*, 206:305-312 (1989); Posfai, et al., *Gene* 74(1):261-265 (1988)). This method is most useful when three conditions obtain: a cognate methyltransferase exists, the genes for the two functions are tightly linked in the DNA, and the methyltransferase is not effectively expressed in *E. coli*. Because the methyltransferase is not effectively expressed, the methyltransferase selection method cannot be used. Briefly, this alternative method is as follows: the polypeptide sequence of the conserved polypeptide motif elements is reverse-translated into a pool of DNA sequences each capable of specifying the polypeptide sequence in question. This pool is called a degenerate pool, because the genetic code is degenerate--several different DNA triplets can specify the same amino acid in many cases. This degenerate pool of oligonucleotides is then used to amplify fragments of DNA from genomic DNA or from a clone library. The sequence of the PCR fragments is then determined, enabling design of further non-degenerate (unique) primers that detect the presence of the proper sequence in the genomic DNA or the clone library by hybridization or PCR. Adjacent DNA sequence can then be obtained by the inverse-PCR method or by Southern blot screening procedures; further sequence can be determined; and finally the complete restriction system can be assembled. This method can be used either alone or in combination with other procedures (below) to isolate the methyltransferase gene and the adjacent endonuclease gene.

"Methylase indicator" DNA damage method.

Another method for identifying clones containing methyltransferase genes (Piekarowicz, et al., *J. Bacteriol.* 173:150-155 (1991); Piekarowicz, et al., *Nucleic Acids Res.*, 19:1831-1835 (1991); Piekarowicz and Weglenska *Acta Microbiol. Po.*, 43(2):229-231 (1994)) relies on methylation-dependent restriction systems McrA, McrBC and Mrr (Heitman and Model, *J. Bacteriol.* 169(7):3243-3250 (1987); Heitman and Model, *Gene* 103:1-9 (1991); Waite-Rees, et al., *J. Bacteriol.*, 173(16):52-7-5219 (1991); Raleigh and Wilson, *Proc. Natl. Acad. Sci. USA* 83:9070-9074 (1986); Kelleher and Raleigh, *J. Bacteriol.*, 173(16):5220-5223 (1991)) and on the *dinD1::lacZ* operon fusion, to enable a method to screen for clones that contain methyltransferase genes. Strains with temperature sensitive mutations in *mcrA*, *mcrBC*, and *mrr* are permissive at high temperature for expression of methyltransferase activity by cloned foreign genes. When these restriction functions are active however (at low temperature), they will cleave DNA methylated by foreign methyltransferase enzymes. This cleavage leads to generation of a signal that induces expression of the endogenous DNA damage inducible (SOS) regulon. The *dinD1::lacZ* transcriptional fusion between one of the genes in this regulon (*dinD*) and the *lacZ* gene is then induced, and  $\beta$ -galactosidase is expressed. Action of the  $\beta$ -galactosidase allows the colonies turn blue on plates containing Xgal. Thus, colonies from a clone library that are white (or light blue) at high temperature but dark blue at low temperature are methyltransferase clone candidates.

N-terminal sequence/degenerate PCR method

It may occur that a methyltransferase gene cannot be identified, or that a methyltransferase gene can be identified but the open reading frame specifying the endonuclease is uncertain. In these cases, an additional useful procedure for identifying the gene for the endonuclease specifically can be applied when the endonuclease can be purified in sufficient quantity and purity from the original organism. In this method, the endonuclease polypeptide is purified to homogeneity and subjected to N-terminal polypeptide sequencing. The polypeptide sequence is reverse-translated into a pool of DNA primers capable of specifying the appropriate sequence, and these primers are used to amplify a

portion of the endonuclease gene from genomic DNA of the original organism or from a clone library.

5 This procedure can be used alone to obtain a portion of an endonuclease gene, or in combination with other methods, such as the degenerate methyltransferase-motif PCR method (Morgan, U.S. Patent No. 5,543,308 (1996)) to obtain portions of genes for both components of the restriction system. The complete genes can be assembled with the assistance of Southern blot or by  
10 further direct or inverse PCR methods. If the cognate methyltransferase gene cannot be obtained or cannot be expressed, the stability and utility of solo endonuclease clones will be severely compromised. Such clones can be stabilized with the use of heterospecific methyltransferase genes, which were not associated with the endonuclease in the original host, if they recognize the same or a related sequence and prevent the endonuclease from cleaving its recognition sequence  
15 (Wilson and Meda, U.S. Patent No. 5,246,845 (1993)).

#### Endo-blue method

Another method for identifying the presence of an endonuclease gene in a clone library, independently of the presence of the cognate methyltransferase gene,  
20 is to introduce the library into a restrictionless host *E. coli* strain containing a reporter of DNA damage. This method is related to "methylase indicator method" above, but the strain used contains no restriction activity specific for methylated DNA. In this case, cleavage occurs due to expression of the restriction enzyme, thereby inducing the SOS regulon (and the *dinD::lacZ* indicator) directly rather than  
25 through the action of the methyltransferase and endogenous restriction activities. Action of the  $\beta$ -galactosidase then allows the colonies to turn blue on plates containing Xgal.

30 This indicator can be used to identify restriction endonuclease clones when a modification methyltransferase gene is poorly expressed, so that some DNA damage occurs despite its presence, or without the methyltransferase when conditional activity of the endonuclease can be obtained. For example, the endonuclease in question may be inactive at low growth temperatures but somewhat active at higher growth temperatures. The latter situation obtains, for  
35 example, with some restriction endonucleases originally expressed in

hyperthermophilic organisms, which normally grow at very high temperatures (Fomenkov, et al., U.S. Patent No. 5,498,535 (1996); Fomenkov, et al., *Nucleic Acids Res.* 22(12:2399-2403 (1994)).

Background of regulation of gene expression in cloned genes.

5 Regulation of expression from vector promoters

In very many instances the problem for the experimenter is to obtain sufficient expression from cloned DNA to enable useful amounts of a gene product to be made in the new cellular environment. Accordingly, there are many expression vectors available that provide one or more promoters enabling high-level transcription activity proceeding through the location at which foreign DNA is to be introduced. Frequently these vectors are provided with a gene for a regulatory molecule such as a repressor of transcription able to regulate expression from the promoter provided, or are used in host organisms that themselves provide such a regulator. In this way, the expression desired can be provided on demand, ie. during induction of specific expression. Many such vectors are described in the art (Sambrook, et al., Molecular Cloning: A Laboratory Manual (1989)).

In some instances, the reverse problem occurs: the product expressed from the cloned DNA is toxic to the cell expressing it for some reason, and ordinary vectors designed for expression at high levels express too much of the toxic product, even in the absence of specific induction. Accordingly, vectors have been described that are designed to express cloned genes at extremely low levels in the absence of induction. The best known of these is the T7 RNA polymerase-dependent expression system designed for use in *E. coli* (Studier, et al., *Meth. Enzymol.*, 185:60-89 1990)). In this system, cloned genes are expressed from a promoter of transcription that is not recognized at all by any endogenous *E. coli* RNA polymerase holoenzyme. Instead, the promoter employed is recognized by the RNA polymerase of bacteriophage T7. This polymerase is not encoded in the *E. coli* genome. This system enables the construction of a clone with toxic properties in the absence of the required RNA polymerase. The clone can then be introduced into a suitable strain into which the T7 RNA polymerase gene has been introduced previously, or the polymerase gene can be introduced by infection with a phage-borne clone.

## Inhibition of expression from indigenous promoter-like sequences

An additional problem with toxic proteins can be encountered when the foreign DNA, introduced into the expression vector, itself contains sequences recognized by the *E. coli* expression apparatus. The specific regulators provided by the vector/host combination will not regulate promoter activity originating within the cloned sequence. In some cases this expression may be the result of specific promoter recognition, but it may also arise simply from adventitious promoter-like activity in DNA, particularly in DNA rich in A+T (Miller and Simons, *Mol. Microbiol.*, 4(6):881-893 (1990)). In such instances a useful method of control is to provide, in the vector, a regulatable promoter opposing the direction of translation of the cloned DNA (Cole and Honore, *Mol. Microbiol.* 3(6):715-722 (1989); Adhya and Gottesman, *Cell* 29(3):939-944 (1982); Elledge, et al., *Proc. Natl. Acad. Sci. USA*, 86(10):3689-3693 (1989); Simons, and Kleckner, *Annu. Rev. Genet.*, 22:567-600 (1988); Roberts, et al., International Publication No. WO 99/11821 (1999)). A high level of transcription in the direction opposite that needed for polypeptide expression can interfere with expression in at least two ways. First, it can occlude transcription in the direction needed for expression; and second, it can prevent translation by allowing formation of RNA-RNA hybrids between the RNA used for expression of the toxic protein and the RNA directed in the opposite sense (antisense RNA).

## Cloning into an expression vector for tight regulation

Restriction endonucleases, which cleave DNA at particular sequences, are normally associated with protective modification methyltransferases. In the present method it is quite likely that the gene for such an endonuclease will be isolated without its partner methyltransferase gene. Very tight regulation of the cassettes thus cloned is therefore critical.

A convenient tightly regulated expression plasmid, pLT7K, is available into which pooled PCR fragments can be cloned (Roberts, *supra* (1999)). In this vector, two levels of control are available: expression is inducible and inhibition is repressible. A T7 gene 10 promoter reads into one side of the cloning site; LacI provided by the vector represses expression from this promoter, as is expression

of the T7 RNA polymerase provided by the host cells used for expression. Further control can be obtained by the use of pLysP, which expresses an inhibitor of T7 RNA polymerase.

5 To further reduce expression directed by the cloned fragment, and residual leaky expression from the T7 promoter, the  $\lambda$  pL promoter reads into the other side of the cloning site, antagonizing expression from pT7. This antagonistic transcription is regulated by  $\lambda$  cI<sup>857</sup>, a thermosensitive repressor. At 40°C and in the absence of IPTG therefore, essentially no expression was observed; at 30°C,  
10 some leaky expression is seen; at 30°C in the presence of IPTG, moderate levels of expression can be achieved. This vector has successfully been used to establish the *pacIR* and *nlaIIIR* genes (encoding the restriction enzymes PacI and NlaIII) in the absence of methyltransferase protection, and to express the genes.

## 15 SUMMARY OF THE INVENTION

A general object of the invention is to provide a procedure for obtaining clones of diversity-selected genes. A specific object of the invention is to provide a method for identifying a repeat sequence suitable for identification and cloning of  
20 gene cassettes found in arrays and separated by repeat sequences. A specific example of such a repeat sequence family with 74 members is provided together with the sequences of four contiguous DNA stretches comprising one or more cassette arrays. A further specific object of this invention is to provide a procedure for cloning cassettes from such arrays, by PCR directed by oligonucleotides  
25 hybridizing with the repeated sequences flanking the cassettes. A specific example of such a PCR procedure is provided. A further specific object of this invention is to provide a procedure for cloning such PCR fragments into an expression vector able to stabilize toxic genes such as restriction enzymes. A specific example of such a gene clonable by this procedure is provided. A further specific object of the  
30 invention is to provide a means of identifying particular cloned genes of interest. Accordingly, three methods of identification are provided: one method relies on identification by means of protein sequence similarity; a second method relies on an indirect report of gene activity; a third method relies on direct test of biochemical properties. In accordance with this method, two novel strains that enable provision  
35 of indirect report of expressible cloned nuclease genes in the context of the vector



pLT7K are provided, together with a method of use. A further specific object of the invention is to provide a method for obtaining expression clones of active restriction enzyme genes without prior knowledge of their biochemical activity or DNA sequence. A specific example of a procedure for obtaining such a clone is provided.

Since the invention relates to genes found in a particular sort of hypervariable locus, a description of what sorts of genes these will be is provided.

### Features of gene cassettes useful for cloning methods.

In the particular case of hypervariable loci that are integrons or superintegrons, these regions provide a mechanism for discovery of diversity-selected genes. The features of these systems enable isolation of DNA enriched for certain kinds of genes including restriction enzyme genes, and also enable the cloning, sequencing and expression of products encoded in this DNA.

Three features of cassette arrays are particularly useful for cloning purposes:

- Each gene (rarely, a pair of genes) is embedded in a predictable sequence context--a particular kind of repeated DNA sequence is found on each side.
- Most genes found such arrays are in the same orientation relative to the flanking sequences.
- Expression of cassette-encoded genes is frequently directed from outside the cassette.

These properties make it likely that genes cloned by PCR from the flanking repeat elements will be intact, will be in an orientation specified in advance relative to the cloning vehicle, and can be regulated by expression signals in the cloning vehicle. This yields a set of DNA fragments in which each gene (rarely, a pair of genes) is embedded in a manipulable sequence context--suitable sites for cloning can be included at the 5' ends of the PCR primers.

5 A difficulty with these repeat sequences is that the members of the repeated array are degenerate, so that PCR primers hybridizing to most or all of the members of the array are difficult to design. Accordingly it is important to have available a large number of such sequences, enabling design of multiple family-specific primers. Such a collection of repeat sequences is identified and characterized in accordance with this invention.

10 A second difficulty with these repeat sequences is that individual members of the repeated array display imperfect dyad symmetry elements, making it likely that PCR primers designed will form hairpins or primer dimers and so fail to prime DNA amplification. Accordingly, it is important to design primer that anneal to portions of the repeats that do not display these features. Primers that are able to hybridize with or that enable amplification from many cassettes are provided in  
15 accordance with this invention.

#### **Expression cloning of cassette-encoded genes.**

20 A very large number of uncharacterized cassettes may potentially be obtained by this method, so that the experimenter will require some procedure for sorting through these for functions of interest. Accordingly, the present invention provides a method for obtaining expression of cassette-encoded functions even when toxic, by cloning these into an appropriate vector, such as the pLT7K vector described in International Publication No. WO 99/11821 (Roberts, et al., (1999)).

25 This vector has the advantage (in addition to those provided in the original patent) that it can be used in two configurations in this application. Depending on the orientation of cloning sites on the PCR primers, the expression condition can be either 30 C + IPTG or 40 C - IPTG; and the repressed condition suitably the reverse. This enables flexibility in screening or selecting for molecules that display  
30 activity sensitive to temperature, and in selecting storage conditions for the clone library obtained.

### Strain enabling indirect report of nuclease activity.

A test of function is provided that enables detection of a minority of expression clones of interest in the context of the T7-RNAP dependent regulation required by the vector pLT7K. This test detects nuclease or other DNA damaging activity by SOS induction of *dinD::lacZ* alleles. Two strains are provided:

ER2745: (F<sup>-</sup>  $\lambda$  *fhuA2* [*lon*] [*dcm*] *ompT* *lacZ::T7 gene1 gal* *sulA11*  $\Delta$ (*mcrC-mrr*)114::IS10 *R(mcr-73::miniTn10--TetS)2 R(zgb-210::Tn10 --TetS) endA1*) *dinD2::MudI1734* (KanR, *lacZ*<sup>+</sup>)

ER2746: (F<sup>-</sup>  $\lambda$  *fhuA2 glnV44 e14- rfbD1? relA1? endA1 spoT1? thi-1*  $\Delta$ (*mcrC-mrr*)114::IS10 *lacZ::T7 gene1 dinD2::MudI1734* (KanR, *lacZ(ts)*)

The former can be used at either 30°C or 42°C to indicate DNA damage with a dark blue color against a background of lighter blue colonies. The latter can be used at 30°C up to and including 37°C to indicate DNA damage with blue color of any shade against a background of white colonies. Accordingly, libraries of cassettes cloned into pLT7K (or a derivative) in an orientation such that expression is driven by pT7 in the presence of T7 RNAP and inhibited by expression from  $\lambda$  pL can be screened for activity at 30°C or 37°C (with or without the presence of IPTG) in either strain. Libraries of cassettes cloned into pLT7K in an orientation such that cassette expression is driven by  $\lambda$  pL and inhibited by pT7 can be screened at 37°C (with or without IPTG) in either strain or 40°C (with or without IPTG) in ER2745 but not ER2746. In each case the presence of activity is indicated when a colony turns bluer than the majority class, and when this property is stable upon reisolation as a single-colony derivative of the original transformant.

These strains may similarly be used to indicate DNA damage provoked by any agent, including enzymes that are not nucleases, by chemical agents, or by radiation. These strains are most distinctively useful when the damage produced results pursuant to a regulated change in the state of T7 RNA polymerase expression as provided within these strains.

**Kinds of genes for which this method may be applied.**

5 In accordance with this invention, a limitation is provided for the kinds of  
genes for which the invention is useful. Some kinds of genes are likely to be  
present in cassette arrays, while others are unlikely to be present in them. The  
original cassettes of known function all specified resistance to drugs or other  
antibacterials. There is no a priori reason to suppose that integrons cannot mediate  
10 the spread of functions other than drug resistances. Types of genes likely to be  
enriched in such arrays include functions useful individually or in pairs, and  
subject to highly variable selective value. Typically such genes will be subject to  
strong episodic selection, very important some of the time but not useful at all the  
rest of the time. In some cases they will be episodically essential--necessary for  
cell survival: drug resistance factors, restriction-modification systems. In other  
15 cases they may be episodically of very high selective value, but not necessary for  
survival as such. Examples would include specific adhesins that allow the cell to  
attach to a particular surface in a rich environment; specific enzymes that modify an  
abundant material in the cellular environment to convert it to a form usable as  
nutrition; or specific toxin molecules that interact with a host organism. Many  
20 individual members of a particular species will elaborate gene products that have  
common general properties (adhesins stick to things). An important feature of  
relevant gene products, however, is that among the population will be found  
examples with highly variable specificities (there are many different kinds of  
specific surfaces to stick to, from rocks to intestinal mucosa to urinary epithelium).

25  
Cassette arrays therefore will be enriched for genes that are subject to  
selection for diversity as described above: that is, genes that are advantageous  
when rare but of no particular use when frequent in the population; and those  
episodically required.

30  
Types of genes expected to be absent from such arrays include all of the  
basic components of the cellular maintenance machinery: DNA replicases, basic  
transcription factors such as vegetative RNA polymerase, the translational  
machinery, enzymes of small molecule metabolism central to cellular physiology  
35 such as those of the tricarboxylic acid cycle. They should be absent for two

reasons. First, no selective advantage is expected from maintaining variability as such in the pool of alleles available to a population of cells. Second, many such proteins must maintain (conserve) specific interactions among several different proteins (replicase/RNA polymerase/translation initiation factor interactions for example).

## BRIEF DESCRIPTION OF THE FIGURES

Figure 1 is a schematic of the structure of characterized integrons, arrays of gene cassettes (thin lines; fn1, fn2, fn3) separated by repeated sequences (filled boxes; 59 bp elements). These are assembled by the action of a site-specific integrase (large box; intI) by insertion into attI (arrows) of extrachromosomal circles (cassette). Cassettes are transcribed from a promoter within the integrase gene (arrow). Many integrons are associated with a conserved sulfonamide resistance gene (sulI) that is not part of the integron itself.

Figure 2 is a schematic diagram of a fragment of a superintegron identified in *Vibrio cholerae*. Open reading frames (1-9 and mrhA, mrhB) are separated by repeats (boxes) that are similar to 59 bp elements of integrons

Figure 3A-3E is an alignment of some of the PAR elements (SEQ ID NO: 96 through SEQ ID NO:116), those identified in superintegron contig 1 (SEQ ID NO:1) by the motif search procedure described in Example 1. Consensus lines show bases shared by all (top line), 90% (second line) or the majority (third line) of the elements in the alignment. Individual entries are the same as the majority consensus except for the bp shown.

Figure 4. is a dotplot display illustrating an alternative method for identifying repeated sequences.

Figure 5. illustrates the self-complementarity of an individual PAR element (SQUIGGLE display of the output of FOLD in the GCG program set).

Figure 6 illustrates alignments of subfamilies identifiable in the set of PAR elements herein (SEQ ID NO:5 through SEQ ID NO:78) shown in Table 1. Panels

A-D, families 1-4. Each family alignment includes PAR2 as an outgroup member, since PAR2 is the most distantly related of the elements identified. Families were identified as bushy groups in a phylogenetic tree generated from the CLUSTAL alignment of the 74 elements.

Figure 7 illustrates the location of oligonucleotides used for Southern blots (panel A) and PCR fingerprinting (Panel B) in relation to the majority consensus of all PAR elements and in relation to a typical cassette.

Figure 8 illustrates a Southern blot hybridization of a mixture of Oligonucleotides 2-5 (SEQ ID NO:79 through SEQ ID NO:83; Fig 7, see also Table 2) to *P. alcaligenes* DNA.

Figure 9 displays an agarose gel of PCR products generated from chromosomal DNA of isolates of six *Pseudomonas species* by the use of oligonucleotides 6 and 7 illustrated in Fig. 7.

Figure 10 illustrates the scheme for forming a clone library of cassette-encoded open reading frames and expression of their products from pLT7K.

## DETAILED DESCRIPTION OF THE INVENTION

In accordance with one embodiment of the invention, there is provided a novel method for the direct cloning and expression of diversity-selected genes residing in cassette arrays. In general, the method comprises the following steps, although as the skilled artisan will appreciate, modifications to these steps may be made without adversely affecting the outcome:

**1) The class of genes of interest is identified and the suitability of the class for the method is evaluated.**

In one embodiment of the invention the desirable genes are those for restriction endonucleases and modification methyltransferases. Types of genes likely to be enriched in cassette arrays include functions useful to the organism individually or in pairs, and subject to highly variable selective value. A function may be identified as likely to be encoded by genes in such arrays when a survey of different isolates of a species determines that the presence of the function, or its

specificity, is variable within the collection of isolates. For example, a survey of isolates of *Escherichia coli* reveals that many isolates but not all isolates express type II restriction enzymes; and that of those that do, the specificity of the enzyme (the sequence recognized) is variable, with many different specificities determined within the species. Candidate functions that will be subject to such variation include, in addition to restriction enzymes, cell surface antigens such as polysaccharide antigens or polypeptide antigens or secreted molecules; adhesins of various sorts such as fimbrial proteins, pilus proteins or outer membrane proteins; transporters of small molecules, especially those with narrow specificity; exported functions such as toxins, hemolysins, hemagglutinins, kinases and signalling molecules; detoxifying enzymes such as drug resistance determinants; catabolic enzymes specific for compounds episodically available (excluding those required for central metabolic pathways such as the tricarboxylic acid cycle); enzymes for biosynthesis of rare sugars (excluding those required in all cells, such as ribose, deoxyribose, and sugars of the cell wall), especially of those sugars that form part of the pericellular envelope.

In one embodiment of the invention, the desirable genes are those for restriction endonucleases and modification methyltransferases. Typically such genes will be subject to strong episodic selection, very important some of the time but not useful at all the rest of the time. Restriction functions can provide a very powerful protection against the invasion of foreign DNA (as when a bacteriophage infects the cell). This protection will only be effective if the host from which the bacteriophage did not carry the same restriction functions--otherwise its DNA would already carry the protective modification pattern of the invaded cell. Populations should therefore carry a wide variety of specificities of restriction-modification systems, and should switch them rapidly on an evolutionary time-scale. In accordance with this expectation, many restriction systems are found on plasmids. Integron-like structures provide an easy way to acquire a restriction system from a foreign source such as a plasmid, which might not establish itself successfully. The existence of the repeat elements would also provide a mechanism for a high rate of loss (by unequal crossing-over or slipped-mispairing during replication), thereby conferring a high degree of fluidity upon the cell's complement of restriction-modification systems.

## 2) DNA preparation

Genomic DNA is prepared from a strain of interest or from a consortium of strains or from an environmental source by methods known in the art, or DNA of plasmid, cosmid, BAC or PAC clones of genomic DNA from such sources is prepared.

## 3) Suitability of the DNA preparation for use of the method.

This is evaluated by determining the presence of repeated sequence arrays. Preferred methods are Southern blot hybridization or PCR fingerprinting using hybridization probes or PCR primers listed in Example 1. Other suitable primer pairs may be designed based on sequences listed in Example 1, or on other particular repeat sequences identified by methods described in Example 1. A DNA preparation is suitable for use if a hybridization signal is obtained or PCR products are obtained or both. In a preferred embodiment, PCR conditions are optimized using a non-proofreading DNA polymerase, by varying primer-template ratio, annealing temperature, magnesium ion concentration and extension time.

## 4) Cassette isolation

The DNA preparation is subjected to PCR employing a pair of primers annealing to repeat sequences flanking the cassettes and containing at their 5' ends sites for endonucleases compatible with cloning into a plasmid vector. Preferred primer pairs include those listed in Example 2; other suitable primer pairs may be designed based on sequences listed in Example 1, or based on other particular repeat sequences identified in the literature or by methods described in Example 1. In a preferred method, PCR conditions are optimized using a proofreading DNA polymerase, by varying primer-template ratio, annealing temperature, magnesium ion concentration and extension time. PCR fragments are purified away from primers, for example by means of size fractionation using commercially available kits.



## 5) Cassette cloning

The PCR fragments are digested with the appropriate restriction endonucleases for cloning, in one preferred procedure with XhoI and XbaI. The digested fragments are ligated into a suitable vector. Preferred vectors for this purpose have two particular properties. First, they contain a cloning site disposed to allow directional cloning of fragments. Directional cloning methods include the process of digesting the vector with two different restriction enzymes such that the single-stranded extension at one end does not hybridize the single-stranded extension at the other end of the vector backbone containing the origin of replication; and then ligating, to that vector backbone, DNA fragments having an extension at one end that hybridizes with one single-stranded extension of the vector backbone, and having an extension at the other end that hybridizes with the other single-stranded extension of the vector backbone. Other directional cloning methods can be envisioned, including for example the use of site-specific recombination enzymes, or hybridization of extensions provided by methods other than restriction enzyme cleavage. Second, preferred vectors contain two independently regulatable expression signals, one on each side of the cloning site described above and directed toward expression of the sequence resident at the cloning site. One preferred vector is pLT7K (Roberts, et al., International Publication No. WO 99/11821 (1999)). Other vectors include pBR322, pUC19, pACYC184, pSC101, pBeloBAC11, or their derivatives.

## 6) Strain choice

The ligated products are transformed into a strain suitable for screening or selecting for cassettes encoding desirable functions. For this purpose the strain must be compatible with the expression regulation signals provided by the vector chosen and must enable the method to be used for identifying desired cassettes.

In the simplest case, sequencing large numbers of cloned cassettes and subsequently evaluating the sequence information will identify cassettes of interest by bioinformatic methods. Such methods include matching the cassette-encoded sequences against public or private databases by means of similarity-determining algorithms such as BLAST or FASTA, or by employing a motif or pattern-based

search of the cassette-encoded sequences employing databases such as the PROSITE profiles database or the BLOCKS and PRINTS databases (Patterson, M. and Handel, M. (1998) Trends Guide to Bioinformatics, Elsevier Science, Cambridge, UK). In this case there are few constraints on strain or vector choice.

In other cases, cassettes of interest will be identified by sequence-based methods such as PCR or hybridization with probes. In these cases there are also few constraints on strain or vector choice.

In a preferred embodiment, cassettes of interest will be identified by activity expressed in vivo. In this case the choice of strain and vector is constrained: vector and strain must be compatible, enabling suitable regulation of cassette expression; by the nature of the activity to be expressed will also constrain strain choice.

In one embodiment, the activities to be expressed are modification methyltransferase activity or restriction endonuclease activity, both of which are amenable to identification by indirect report of activity based on damage inflicted in intracellular DNA and induction of the DNA damage repair response. Two preferred strains ER2745 ( $F^- \lambda$  *fhuA2* [*lon*] [*dcm*] *ompT* *lacZ::T7 gene1 gal* *sulA11*  $\Delta$ (*mcrC-mrr*)114::*IS10* *R*(*mcr-73::miniTn10--TetS*)2 *R*(*zgb-210::Tn10 --TetS*) *endA1*) *dinD2::MudI*1734 (*KanR*, *lacZ*<sup>+</sup>), and ER2746: ( $F^- \lambda$  *fhuA2 glnV44 e14- rfbD1? relA1? endA1 spoT1? thi-1*  $\Delta$ (*mcrC-mrr*)114::*IS10 lacZ::T7 gene1 dinD2::MudI*1734 (*KanR*, *lacZ*(*ts*)) are strains compatible with the vector pLT7K.

ER2745 is derived from the particular strain background normally used for T7 RNAP-directed expression, and is ultimately a derivative of *E. coli* B. The protein expression properties of this strain background are well understood. This strain is transformable with DNA, but the level of transformation obtained is less than with other strains. The amount of the indicator *lacZ* expressed in the absence of DNA damage is relatively high, leading to light-blue colonies on Xgal plates even when no damage has occurred.

ER2746 carries a thermosensitive *lacZ* moiety. This is useful because it lowers the light-blue background color observed on X-gal by the original *dinD* indicator allele. Discrimination between clones inducing some damage (which are

of interest) and those inducing no damage (which are not) is improved in this situation. However, this allele cannot be used to detect DNA damage at high temperature ( $>37^{\circ}\text{C}$ ), because the *lacZ* moiety of the indicator fusion is inactive, and will remain white even in the presence of extensive DNA damage. This was demonstrated by testing at various temperatures for induction of blue color by nalidixic acid, a well-characterized DNA damaging agent, on plates containing X-gal.

Further refinement of this system is possible; for example, transcriptional fusion of a drug-resistance gene to a damage-inducible promoter should allow selective isolation of clones of interest, rather than the more-laborious screening procedure. Use of a variety of drug concentrations would then allow isolation of clones with different levels of DNA-damaging activity. Introduction of a *recD* mutation would inactivate the major ATP-dependent double-strand exonuclease of the cell, while an *xth* mutation would inactivate ExoIII, the major ATP-independent double-strand exonuclease. A triply nuclease-deficient strain should be viable but may not stably maintain the plasmid (Niki, et al., *Mol. Gen. Genet.* 224(1):1-9 (1990)).

Other DNA damage-inducing promoters that can be used include those identified by (Lewis, et al., *J. Bacteriol.*, 174:3377-3385 (1992); Lewis, *J. Mol. Biol.*, 241:506-523 (1994)): these are promoters of *recA*, *lexA*, *uvrA*, *uvrB*, *dinG*, *polB*, *uvrD*, *ruvAB*, *umuDC*, *sulA*, *dinH*, *dinI*, *sosA*, *sosB*, *sosC*, *sosD*. Other SOS-inducible genes identified include *recN*, *dinB* and *dinF* (Walker, *Microbiological Review*, 48:60-93 (1984)). Some other indicator/reporter genes that can be used were reported in (Fomenkov, et al., *supra* (1995)).

#### 7) Cassette identification: endonuclease genes

Following transformation or electroporation of the cassettes ligated with the chosen vector into the chosen strain, transformants are plated onto suitable media. In the preferred procedure, the vector is pLT7K, the strain is ER2746, plates are Luria-Bertani plates with ampicillin, and incubation is at  $40^{\circ}\text{C}$ . Colonies are replica plated onto plates containing Xgal with or without IPTG (at concentrations varying from 0.1 mM to 1 mM) and one set of replicas is incubated at each of three

temperatures, 30°C, 37°C and 40°C. These conditions range from fully inducing and indication-capable (30°C, high IPTG) to fully repressing and indication-negative (even induced cells would not turn blue due to the thermosensitive *lacZ* allele) (40°C, no IPTG) Colonies that are blue at any condition are then candidate  
5 nuclease genes. The darker the blue color, the greater the DNA-damaging activity.

Individual colonies can then be recovered from master plates that have not been subjected to the damaging condition, to assure recovery of the original sequence, grown in small cultures (10 ml LB with antibiotic) and plasmid  
10 preparations made for storage.

Reversing the configuration of expression so that the repressing condition is at 30°C +IPTG and the inducing condition is 40°C - IPTG can be easily accomplished with pLT7K by switching the cloning sites added to the  
15 oligonucleotide primers for PCR so that cassettes are in the reverse orientation. This may be desirable to facilitate storage of never-induced colonies. For this purpose strain ER2745 is the preferred strain, since the damage-inducible fusion carries a wild type *lacZ* allele that enables indication at 40°C. In that case, the colonies desired will be darker blue than the normal light blue color.

Further characterization is then carried out on the identified plasmids, either continuing from the replica plate masters or from the archived plasmid DNA following retransformation. Further characterization includes some or all of the  
20 following three steps.

Crude extract assay: Clones positive in the DNA-damage screen are grown at in medium-sized cultures (20-200 ml) at 40°C -IPTG (noninducing conditions) in LB + ampicillin to late log phase, and shifted to the inducing condition identified for the clone (usually 30°C + IPTG, but possibly a semi-inducing condition) for four hours. This procedure was successful in allowing expression of an amount of PacI similar to that expressed in the native host, *P. alcaligenes* (D. Byrd, personal communication). Cells are then collected by centrifugation, resuspended in buffer, lysed by lysozyme-EDTA treatment, and clarified by centrifugation.

Crude extracts supernatants are then assayed for nuclease activity in a general screen for 4-6 base cutters, using standard plasmid, phage and viral DNAs such as pUC19, pACYC187, pACYC177, pBR322, M13mp18 replicative form DNA, lambda DNA or T7 DNA at 37- 68 °C. Some 8-base specificities may be detected by this method as well.

DNA digestion patterns are resolved by agarose gel electrophoresis using an agarose concentration suitable for visualization of bands between 200 and 0.05 kb (usually 0.7% agarose and 1.3 % agarose), and detected by ethidium bromide staining.

DNA digestion patterns are then evaluated and the recognition sequence is determined by methods known in the art. Further purification of the endonuclease thus identified may be required for these methods to be applied.

Crude extract supernatants are also assayed in an in vitro screen for enzymes with 8-base sites, using chromosomal DNAs of varying GC-content: *Rhodobacter sphaeroides*, *Escherichia coli* and *Staphylococcus aureus* range from 66% to 34% G+C and are suitable for detecting a variety of enzymes with rare sites. It is usually possible to distinguish between nonspecific nuclease and an 8-base endonuclease, since specific fragments (especially large ones) are not subject to further digestion; even though the fragments are not resolvable on the gel (and the recognition site cannot be deduced), the result is recognizably different from that produced by nonspecific nucleases (which preferentially degrade large fragments). In each case, aliquots of extract are incubated with potential DNA

substrates in the presence of  $Mg^{++}$  and resolved on agarose gels followed by ethidium bromide staining.

Isolates that yield a positive result on chromosomal digests but not in digests of standard substrates are then further characterized by searching for alternative substrates, guided by the G+C content of the chromosomal DNA yielding a positive result.

Pulsed-field gel assay: A potentially more-informative assay for 8-base recognition sites relies on separation of total chromosomal fragments on pulsed-field gels. When crude extracts are used for screening procedures, these gels are too cumbersome and too sensitive to other nucleases in the extract to be generally useful.

In standard procedures, the substrate DNA is obtained by first embedding whole cells in agarose plugs. DNA is released from the cells in situ by means of a series of enzymatic treatments and washes that degrade the cell wall. The restriction endonuclease is then incubated with the plug; this usually takes several hours, since the enzyme must permeate the agarose and the remnants of the previous digestions.

In this method the restriction nuclease digestion step consists of inducing expression within the cell, before agarose is added; embedding the cells in agarose and subjecting the cells to electrophoresis on a pulsed-field agarose gel. Controls include: positive control, standard digestion of the host DNA embedded in agarose plugs with purified *PacI* and *NotI*; and negative control, samples of the host containing the empty vector, treated in parallel with the experimental samples.

Possible improvements in the strain used for this part of the survey include introduction of a *recD* mutation, which would inactivate the major ATP-dependent double-strand exonuclease of the cell; and introduction of an *xth* mutation that would inactivate the major ATP-independent double-strand exonuclease. A triply nuclease-deficient strain (*endA xth recD*) should be viable but may not stably maintain the plasmid (Niki, *supra* (1990)).

Isolates identified by this method are then carried further, with further purification and overexpression of the cassette-encoded polypeptide, so that conventional pulsed-field analysis can be carried out.

Fingerprinting: Plasmid DNAs prepared from candidate clones obtained by the indirect report assay are fingerprinted by restriction enzyme digestion. Each candidate is digested separately with two to four enzymes with four-base recognition sites: in the preferred example, with HaeIII and MseI to yield a patterns characteristic of the cloned cassette.

Sequencing: All plasmids that result in banding patterns in crude extract or pulsed-field gel assays are then sequenced.

All fingerprinted plasmids are grouped according to fingerprint and two in each class are sequenced. A minimum of three-fold sequence coverage will be required in order to have sufficient confidence to carry out preliminary homology searches.

Sequencing is carried out using the Tn7-based transposition system, GPS<sup>TM</sup>-1 (NEB Catalog No. 1700, New England Biolabs, Inc., Beverly, MA). This system enables introduction of primer-binding sites at random locations in plasmids of interest, rapid mapping of the location of the insertion by digestion with rare-cutters that cleave within the transposon, and sequencing of the insertions within the fragment of interest. With these target molecules, about 20% of transposon insertions will be found within the sequence of interest. No more than 6 suitable insertions are needed in most cases, since cassettes are normally smaller than 2 kb. Two sequence runs (500 bp per run) from flanking vector primers and 12 runs from insertions will yield 7000 bp of raw sequence, approximately 3-fold redundancy. This will be sufficient for primary analysis. Further sequencing can be carried out to obtain high-quality sequence of the most interesting fragments.

Alternative sequencing methods may be used, such as primer-walking, nested deletion construction, or alternative transposon-based methods such as Primer Islands (Perkin-Elmer).

Sequence Evaluation: Homology to genes in public databases will help to exclude candidates for new type II RM genes. Many genes that might be recovered during this procedure exhibit conserved amino acid sequence segments: topoisomerases, helicases, nicking enzymes associated with conjugal plasmid transfer, and transposases all can be found annotated in databases, identified by BLAST or other homology search procedures. Genes for type II restriction enzymes, on the other hand, rarely can be identified in this way. When they can be identified by homology, they are almost always isoschizomers of (recognize the same site as) the enzyme in the database (R. Roberts, personal communication). Thus, the target genes (endonucleases recognizing new specificities) can be expected among those not identified by homology search.

## **2. Cassette identification: methyltransferase gene acquisition.**

In one preferred procedure, the desirable function is a methyltransferase gene, which may be selected or screened for by methods known in the art, described above.

### **A. The methylase selection method**

This may be used if an endonuclease with suitable specificity is available. This method will be applicable when something is known or suspected about the specificity of potential methyltransferase enzymes and a suitable endonuclease is available. Such an endonuclease may be a heterologous endonuclease recognizing a subset of the relevant sites.

### **B. The methyltransferase indicator method**

This may be used if the vector employed is compatible with the strains previously described (Piekarowicz, et al., *supra* (1991); Piekarowicz, et al., *supra* (1991); Piekarowicz and Weglenska *supra* (1994)), with the proviso that the *dinD::lacZ* indicator allele resident in the strains identified in (Piekarowicz and Weglenska, *supra* (1994)) are unable to indicate at temperatures above 37°C, so



only the presence of blue color at or below that temperature should be evaluated. Other strains derived from these may be constructed to enable use of other vectors such as pLT7K.

C. Degenerate methyltransferase-motif PCR

The method of may be employed alone, or the degenerate methyltransferase-motif primers may be combined with a repeat-specific primer or primers annealing to the flanking repeats in a single orientation, such as those employed in PCR fingerprinting or cassette cloning as described above.

D. Biochemical methods

Other methods for evaluating the presence of methyltransferase genes include detection of enzymatic activity such as evaluation of  $^3\text{H}$ -SAM incorporation into specific DNA sequences and may be applied to individual clones or pools of clones.

E. Hybridization methods

Hybridization detection methods such as colony lifts may be employed to detect the presence of genes with high levels of DNA homology to available methyltransferase genes or to oligonucleotides designed based on the sequences of those genes.

The present invention is further illustrated by the following Examples. These Examples are provided to aid in the understanding of the invention and are not construed as a limitation thereof.

The references cited above and below are herein incorporated by reference.

## EXAMPLE 1

## IDENTIFYING REPEAT SEQUENCES AND OBTAINING CASSETTES

5 This Example outlines the general strategy for identifying a candidate repeated sequence. It also provides a specific repeated sequence family, probes for identification of organisms containing similar repeats and primers for amplification of the gene cassettes.

**A) Cloning of portions of a superintegron array.**

10 The organisms expressing PacI and PmeI were isolated by at NEB (Polisson, U.S. Patent No. 5,098,839 (1992); Morgan and Zhou, U.S. Patent No. 5,196,330 (1993)). These restriction enzymes are made by particular isolates of *Pseudomonas alcaligenes* (ATCC No. 55044) (NEB Deposit No. 585, New England Biolabs, Inc.; Beverly, MA) and *Pseudomonas mendocina* (ATCC No. 55181) (NEB Deposit No. 698, New England Biolabs, Inc., Beverly, MA) respectively. The genes encoding these enzymes were identified and cloned using seven steps: 1) PacI and PmeI were purified to homogeneity from *Pseudomonas alcaligenes* (ATCC No. 55044) (NEB Deposit No. 585, New England Biolabs, Inc.; Beverly, MA) and *Pseudomonas mendocina*. (ATCC No. 55181) (NEB Deposit No. 698, New England Biolabs, Inc., Beverly, MA) by the methods of (Polisson, *supra* (1992); Morgan and Zhou, *supra* (1993)). 2) The N-terminal sequences of these proteins were obtained by standard microsequencing methods. 3) Degenerate oligonucleotides, designed on the basis of these sequences, were used to obtain PCR fragments encoding these N-termini. 4) The DNA sequence specifying these N-termini was determined from the PCR fragments. 5) Unique oligonucleotides designed from these specific sequences were used for inverse PCR, to obtain larger fragments encoding the entire genes. 6) In both cases, suitable enzymatic activities were identified in crude extracts of *E. coli* carrying the relevant genes under the control of the T7 RNA polymerase. 7) Further cloning of adjacent sequence was carried out, and sequence was obtained of 4.07 kb of *Pseudomonas alcaligenes* ((ATCC No. 55044) (NEB Deposit No. 585, New England Biolabs, Inc.; Beverly, MA) DNA and 5.37 kb of *Pseudomonas*

*mendocina* (ATCC No. 55181) (NEB Deposit No. 698, New England Biolabs, Inc., Beverly, MA) DNA.

Examination of these sequences by visual inspection enabled preliminary identification of repetitive sequences common to both gene segments. Further cloning experiments were aimed at obtaining a complete sequence description of the cassette array residing in *Pseudomonas alcaligenes* (ATCC No. 55044) (NEB Deposit No. 585, New England Biolabs, Inc., Beverly, MA), resulting in four segments of contiguous sequence as described below. Routine cloning procedures were from (Sambrook *supra* (1989); Maniatis, et al., Molecular Cloning: A Laboratory Manual, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY (1982); Raleigh, et al., Current Protocols in Molecular Biology John Wiley and Sons, New York, pp. 1.4.1-1.4.7 (1989); Moore, et al., Current Protocols in Molecular Biology, John Wiley and Sons, New York, pp 2.0.1-2.6.12 (1999)).

In the expectation that repetitive arrays might be unstable in *E. coli*, we initially avoided attempting to isolate large fragments containing PAR elements. Further *P. alcaligenes* (ATCC 55044) (NEB Deposit No. 585, New England Biolabs, Inc., Beverly, MA) chromosomal DNA fragments were obtained from HindIII libraries constructed by cloning size-selected HindIII fragments into the HindIII site of pBR322. Chromosomal DNA of *P. alcaligenes* (ATCC No. 55044) (NEB Deposit No. 585, New England Biolabs, Inc., Beverly, MA) prepared by the procedure described in the manual of Qiagen (Genomic tip 100/G (Cat 10243) was digested with HindIII to completion. HindIII fragments were isolated by gel fractionation on agarose gels (0.7%) and fragments between 2 kb and 10 kb were isolated using QIAquick Gel extraction kit (Cat # 28704) according to the instructions of the manufacturer and ligated with HindIII-digested dephosphorylated pBR322.

The rationale for this procedure is that *P. alcaligenes* DNA is GC rich while the HindIII site is AT rich (AAGCTT). Therefore few chromosomal DNA fragments are as small (2 kb and 8 kb) as those identified by Southern blot to *pacIR* and PAR-specific probes (see section C1 for this procedure). Plasmid preparations were made from 108 of the colonies obtained following transformation using QIAprep Spin Miniprep Kit Cat #27106. 95 of 108 HindIII

clones (88% ) carried inserts. These were digested with AclI (AACGTT), which cuts within the PAR sequence identified by eye but rarely in the GC-rich *P. alcaligenes* chromosome, and clones were identified that carried exceptionally large numbers of AclI sites. 11% of clones with inserts (11 clones) fit this criterion. Further characterization by PAR-specific PCR (see Section C2) and sequence analysis (below) verified that these did indeed contain PAR sequences.

The high frequency of PAR-containing fragments in the absence of any selection except for size presumably reflects a higher density of HindIII sites within the PAR-containing region than in the chromosome as a whole. We estimate that size selection eliminated about 90% of all chromosomal sequences. If the total genome is 6-8 Mb (Rodley, et al., *Mol. Microbiol.*, 17(1):57-67 (1995); Dewar, et al., *Microb. Comp. Genomics* 3(2):105-117 (1998)) and 10% of this is represented in the size fraction chosen (600-800 kb total), then 100 inserts of average size ~8 kb would be required to cover all of this fraction. A library of this size would of course not contain all fragments exactly once and not all fragments in the fraction are 8 kb. Nevertheless, the incidence of PAR-containing fragments in the library is consistent with the estimated size of the putative superintegron ( $\geq 60$  kb; 10% of 800 kb would be 80 kb).

Additional clones were isolated in subsequent libraries made by digestion with ClaI and cloning into the ClaI site of pBR322. At this stage instability of large fragments did not appear to be a problem, so the DNA was not fractionated but was cloned directly. PAR-positive clones were identified by PAR fingerprinting by the method described in Section C2.

Candidate PAR-containing clones were sequenced with an ABI377 sequencer using dye terminators. Template generation was by a combination method. In a semi-random phase, a Tn7-based transposon (an early version of the NEB GPST<sup>TM</sup>-1 kit, (New England Biolabs, Inc., Beverly, MA, NEB Catalog No. 7100) was used for insertional mutagenesis of clones, and selected insertions were sequenced using universal primers (PrimerN and PrimerS, (New England Biolabs, Inc., Beverly, MA, NEB Catalog No. OS1266 and NEB Catalog No. 1267) designed to sequence from the transposon.. Sequencing was facilitated by limited mapping of insertions, employing rare-cut sites within the transposon. Vector-

insert junctions of primary clones and of a few deletion derivatives were also sequenced using primers annealing to pBR322 (New England Biolabs, Inc., Beverly, MA, NEB Catalog No. 1204 and NEB Catalog No. 1205).

5 This resulted in four sequence contigs totaling 59.4 kb, containing 74 examples of the repetitive sequence. These sequences are SEQ ID NO:1, SEQ ID NO:2, SEQ ID NO:3, and SEQ ID NO:4.

### 10 B) Formulation of a repeated sequence candidate.

15 The specific repeated sequences that are likely to signal the presence of a cassette array can be identified by similarity to those found in known arrays such as the VCR elements of *Vibrio cholerae*, or by computer-assisted analysis of existing sequence information. These sequences were identified by the following procedure, employing computerized search procedures (both UWGCG SEQED and DNASTAR EDITSEQ programs are suitable): the 5' end of the repeat was found by searching for the sequence TAACWA; the 3' end of the repeats were found by searching for the sequence CGTTRR; and the additional constraint was imposed that the 5' base of the 5' element should be not more than 200 bp from the 3' end of the 3' element. This strategy identified 18 repeated elements in this contiguous stretch of 14.144 kb. For comparison, a similar search employing the motifs suggested by Hall (5) identified 11 elements; 10 of these were congruent with the set identified by the strategy cited here, and one aligned very poorly in the internal regions with the others identified by either strategy.

25 Fig 3 shows an alignment of a set of such sequences identified in a part of the *P. alcaligenes* (ATCC No. 55044) (New England Biolabs, Inc. Beverly, MA, NEB Catalog No. 585) superintegron sequence SEQ ID NO:1. The elements were aligned using the DNASTAR MEGALIGN program, by the CLUSTAL method. The alignment shows a majority consensus (third line), a 90% consensus, at which 30 16 of the 18 elements are identical (second line) and an identity consensus, with which all elements agree. Only those positions that disagree with the majority consensus are shown on the alignment. 48% (42/87) of positions in the alignment are identical in 90% of representatives; the most divergent representative (PARf9) still agrees with the majority at more than half of positions (47/87).

An additional method for identifying such a repeat is to use a computerized comparison algorithm such as UWGCG COMPARE and DOTPLOT, or the DNASTAR algorithm ALIGN with the DOTPLOT subprocedure. The output of these programs will identify off-diagonal similar sequences (Fig 4; window of 30, match of 24), which can then be examined more closely using a program feature (in DNASTAR) or by noting the approximate positions of the alignment and following with the UWGCG BESTFIT algorithm on the local subsequences surrounding the diagonal. The DOTPLOT method identified 18 elements also: 16 of these were identified by the strategy cited here while two of those identified by the motif search were not found by DOTPLOT. More sophisticated computerized search procedures based on these methods may also be developed and employed for this purpose.

A complete set of the elements identified by searching for the motifs as described is displayed listed herein (SEQ ID NO:5 through SEQ ID NO:78 Table 1). In these elements, an additional two bp adjacent at the 5' end have been added to each element, since these bp are conserved in the majority of the sequences, as 5' GC 3'. One additional base has been added at the 3' end, since this bp is also conserved as C in the majority of sequences. The length of each element, and its location in the relevant contig, and the name of the contig in which it is found is also entered in this table.

It may be noted that the individual sequences within the set display imperfect internal inverted repetition (Fig 5 shows an example of potential secondary structure). This property was also observed in "59 bp elements" and VCR elements.

It may also be noted that the PAR elements fall into families of more-closely related sequences. Alignments of four of these families are displayed in Fig. 6A-6D. Knowledge of these families will inform the design of specific oligonucleotides for further procedures such as those employed below.

Once a repeat sequence candidate or family has been chosen, either from among known arrays or by analysis of new sequence, oligonucleotide probes and

primers can be designed for use in Southern blot and PCR experiments, described further below. Examples of these are shown aligned with the consensus of 74 PAR elements (majority rule) in Fig. 7A (Oligonucleotides 1-5 (SEQ ID NO:79 through SEQ ID NO:83; see Table 2) for Southern blot) and 7B (Oligonucleotides 6 and 7 (SEQ ID NO:84 and SEQ ID NO:85; see Table 2) for PCR).

### C) Identifying candidate prokaryotic populations.

With the information obtained from one or more array sets, it then becomes possible to screen additional isolates for the presence of such arrays by Southern blot procedures or by PCR.

C1) Southern blot to *Pseudomonas alcaligenes* (ATCC No. 55044) (NEB Deposit No. 585, New England Biolabs, Inc., Beverly, MA)

A Southern blot (Fig. 8) was carried out using a mixture of biotin-labeled oligonucleotides (Oligonucleotides 2-5, SEQ ID NO:80 through SEQ ID NO:83; see Table 2) as a probe for repeat sequences (PAR elements), and chromosomal DNA of *P. alcaligenes* (ATCC 55044) (New England Biolabs, Inc., Beverly, MA, NEB Catalog No. 585) prepared by the procedure of Qiagen (Genomic tip 100/G (Cat 10243). Restriction digests with 8 different restriction enzymes (SphI, PstI, StuI, NdeI, NcoI, EcoRI, ClaI and HindIII) were carried out according to the manufacturer's instructions (New England Biolabs, Inc., Beverly, MA). Products were subjected to electrophoresis for 1 h at 100 mA in 0.7% agarose with Tris Borate buffer (composition 0.09 M Tris-borate, 0.002 M EDTA,  $10^{-4}$   $\mu$ g/ml ethidium bromide). The Southern procedure was carried out according to instructions in the NEBlot® Phototope® kit (New England Biolabs, Inc., Beverly, MA, NEB Catalog No. 7550) using Immobilon-S (Millipore cat #MBBU IMS02) membrane, hybridization at 68°C for 4 h, with 2 washes with at 23°C followed by 2 washes with 0.1XSSPE, 0.1% SDS at 68°C for 5 min. Development was with Phototope®-Star detection kit (New England Biolabs, Inc., Beverly, MA, NEB Catalog No. 7020) chemiluminescent detection according to the manufacturer's recommendations. Fig 8 reveals that multiple fragments in each digest hybridized with the probe, confirming that the oligonucleotide recognized a repeated sequence.

The minimum sum of sizes of hybridizing bands ranged from ~20 (PstI) to ~44 (NdeI) kb, suggesting that a large number of cassettes might be present. Some of these bands may represent doublet or triplet co-migrating species, so the maximum size cannot be reliably estimated.

Alternative possible oligonucleotide sequences might be designed based on specific families of PAR elements. A single oligonucleotide such as Oligonucleotide 1 (SEQ ID NO:79; see Table 2) may be used (data not shown), which may be used to prepare a biotin-labeled probe by starting with an unlabeled oligonucleotide, and labeling it by use of a random-priming kit such as NEBlot® Phototope® kit.

Other detailed procedures may be used for detecting the presence of hybridization between the probe oligonucleotide and the DNA preparation. The Southern blot procedure separates DNA fragments by size, transfers these to a membrane support, denatures the DNA, hybridizes the probe, then separates the hybridized product from the nonhybridized probe (in this case oligonucleotides) by washing. Alternative derived methods for detecting the presence of hybridized DNA include use of arrays of DNA preparations, not separated by size, adsorbed a membrane (dot blots or slot blots (Moore, *supra* (1999)) or microtiter plate (Chaplin and Brownstein Current Protocols in Molecular Biology John Wiley and Sons, New York, Vol. 1, pp. 6.9.1-6.9.7 (1999)) or other support, followed by washing away the unhybridized probe. The configuration of label can be reversed (the target DNA preparation is labeled while the test probe is fixed to the membrane or other support).

Alternative possible detection methods include the use of radiolabeled oligonucleotides (labeled with S<sup>35</sup> or P<sup>32</sup> or P<sup>33</sup>), or of alternative chemical detection methods, such as digoxigenin-based (Roche Molecular Biochemicals Cat #12102201) or fluorescein-based (AP Biotech Cat # RPN 3030) label and detection procedures. Alternative methods of DNA preparation could include purification by detergent/protease treatment followed by precipitation or CsCl centrifugation, or by purification from agarose gels (Moore, *supra* (1999)). Other commercially available kits that rely on gel filtration may also be employed (e.g.



those supplied by 5Prime-->3Prime, or Promega Wizard Genomic DNA Purification Kit, Cat#A1120).

C2) PCR fingerprinting of six *Pseudomonas* species.

5 A second method for detecting cassette arrays in a population is to employ primers annealing to each end of the repeats separating the cassettes in a PCR experiment (Fig 7B and Fig 9). If the repeats are present and close enough to each other for PCR amplification to be effective, DNA bands representing the cassettes will be observed in ethidium-bromide stained agarose gels following  
10 electrophoretic separation.

To validate this method, six species of *Pseudomonas* were tested: *P. maltophilia* NEB Deposit No. 515 (New England Biolabs, Inc., Beverly, MA) (PmlI), *P. fluorescens* NEB Deposit No. 375 (New England Biolabs, Inc., Beverly, MA) (PflMI), *P. putida* NEB Deposit No. 372 (New England Biolabs, Inc., Beverly, MA) (PpuMI), *P. lemoignei* NEB Deposit No. 418 (New England Biolabs, Inc., Beverly, MA) (PleI), *P. mendocina* (ATCC No. 55181) (New England Biolabs, Inc., Beverly, MA, NEB Deposit No. 698), (PmeI) and *P. alcaligenes* (ATCC No. 55044) (New England Biolabs, Inc., Beverly, MA, NEB Deposit No. 585) (PacI). Chromosomal DNA made as above (part A) was used in PCR reactions primed by Oligonucleotides 6 and 7 (Fig. 7; SEQ ID NO:84 and SEQ ID NO:85; see Table 2). PCR reactions included 100 ng DNA, 0.2  $\mu$ mol each oligonucleotide, 1 units of Vent® Exo<sup>+</sup> polymerase, 1X NEB Thermopol buffer in a reaction volume of 50  $\mu$ l. Thermal cycling parameters were 15 sec  
25 denaturation at 95°C, 1 min annealing at 55°C, 1 min extension time at 72°C. 25 cycles were carried out. Products were subjected to electrophoresis for 1 h at 100 mA in 0.7 % agarose with 10<sup>-4</sup>  $\mu$ g/ml ethidium bromide.

30 Figure 8 reveals that two of the six species yielded multiple amplification products from this procedure. This confirms the presence of the repeat segments in the correct orientation and at the correct spacing for amplification to occur. It is not possible to assess the number of potential cassettes from this procedure, since some cassettes may be too long to amplify efficiently, especially in the presence of

shorter cassettes that would be amplified preferentially. In addition, some amplification products may represent amplification across two cassettes. In this case, the repeat separating them might be more distantly related to the primers than those at the ends of the amplicon.

Use of a variety of extension times will facilitate acquisition of a maximum variety of cassette products. Multiple reactions employing alternative primer sets annealing at high efficiency to alternative families of repeats will also increase the total yield of cassettes. Primers 8-11 (SEQ ID NO:86 through SEQ ID NO:89; see Table 2) are candidate primers for the forward direction, while primers 12 and 13 (SEQ ID NO:90 and SEQ ID NO:91; see Table 2) are candidate primers for the reverse direction as displayed in Fig. 8

Alternative methods of visualization include chemiluminescent detection of affinity-labeled oligonucleotide primers, fluorescent detection of fluorescently labeled nucleotides or oligonucleotide primers incorporated during PCR, or autoradiography when using radiolabeled oligonucleotide primers or radiolabeled dNTP.

### C3) PCR fingerprinting of mixed populations

In principle, it should be possible to apply the PCR-fingerprinting strategy to mixed populations to identify the presence of cassette arrays in a minority of the population. At least two kinds of applications to mixed populations can be tried: PCR using combinatorial pools of individual strains, and PCR using environmental DNA.

#### C3a) PCR on combinatorial pools:

Combinatorial pools can be achieved by arraying individual strains in addressable arrays, for example, 96-well plates. Pools can be made combining the individual strains, e.g. all strains in one row in one pool; or all strains in one column in one pool; or all strains in one 2D address from a series of plates. Many such pooling procedures have been worked out and will be familiar to one skilled in the art (e.g. (Chaplin and Brownstein, *supra* (1999); Green, et al., Cloning

Systems, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY, Vol. 3, pp. 297-548 (1999)).

DNA can be made from these strains individually and the DNA samples then pooled; or the strain cultures can be pooled and DNA made from the pool. Each procedure has disadvantages; in the first instance, a larger number of DNA preps must be made; but in the second procedure, different strains may be differentially subject to cell breakage and DNA extraction, and therefore DNA from some strains will be under-represented relative to others.

In such a pooling procedure, some simple controls will allow assessment of the effectiveness of the overall procedure. For example, a positive control--a strain known to contain an array (such as *P. alcaligenes* (ATCC 55044) (NEB Deposit No. 585, New England Biolabs, Inc., Beverly, MA)--can be included in one pool as a single member while the other members are drawn from negative controls--strains known not to contain a responsive array (such as *P. lemoignei* (NEB Deposit No. 418, New England Biolabs, Inc., Beverly, MA). In another, the positive control can be included in duplicate, in another in triplicate, with corresponding reduction in the representation of the negative control. This will enable assessment of the sensitivity of the overall procedure.

#### C3b) PCR on environmental samples:

A DNA source of great interest is likely to be DNA isolated from environmental samples (e.g. soil, water, filtered air etc) without first obtaining organisms in pure culture. In this case, PCR from cassette arrays may be even more desirable as a mechanism for obtaining genes in intact form. In this case, the same kinds of positive and negative controls as those described in C1 may be included. In addition to a dilution series of the positive control in a known negative control, other controls should be included. The original environmental sample from which DNA is to be isolated can be divided and a portion doped with a small amount of the positive control strain. DNA extraction from the sample will then include some of the positive control, enabling that portion of the sample to be used as a control for the efficiency of DNA extraction and recovery of known cassettes from a known source. Inclusion of a dilution series of purified positive

control DNA in the environmental sample DNA will serve as a control for inhibitory materials in the environmental sample.

5 An additional series of controls can estimate the fraction of the sample that derived from eukaryotic organisms. PCR controls can test for the presence of mitochondria, chloroplasts, and nuclear ribosomal DNA genes by methods known to those skilled in the art (von Wintzingerode, et al., *FEMS Microbiol. Rev.* 21(3):213-229 (1997); Sekiguchi, et al., *Microbiology*, 144 (Pt. 9), 2655-2665 (1998)).

#### 10 D) Cloning the DNA fragments.

15 Once DNA fragments flanked by repeat segments have been obtained, these can be cloned by standard methods. PCR products can be purified using the QIAquick PCR purification kit (Qiagen Cat No. 28104) or other similar kits. Fragments can be digested to provide ligatable ends compatible with appropriately-digested plasmid or bacteriophage vectors. In the present Example, XhoI and XbaI sites added to the 5' ends of the oligonucleotide primers used for PCR provides directional cloning into pLT7K (Example 2 below) such that a defined orientation is obtained relative to vector-borne expression signals. Accordingly, the use of regulatory signals residing in the vector is feasible. If regulation of expression is not a concern, any vector can be used to clone such cassettes, provided that suitable cloning sites are included at the 5' ends of oligonucleotides used for PCR. Such vectors may be high-copy (e.g. pUC19), intermediate-copy (e.g. pACYC184 or pBR322), or low-copy (e.g. pBeloBAC11) plasmid 25 replicons, or may be bacteriophage replicons (e.g.  $\lambda$ gt11). Such vectors may contain expression signals suitable for regulated expression in *E. coli* (e.g. pLT7K; see Example 2), or may be designed for expression in an organism suitable for further experimental test of a particular cassette (e.g. *Bacillus subtilis*, *Streptomyces coelicolor*, *Agrobacterium tumefaciens* or other prokaryotic organism).

30 The ligated fragment pool will normally be recovered as a clone library of fragments consisting of colonies of the recipient organism containing one or more

selectable marker of the vector on solid media following transformation by chemical methods or by electroporation (Hanahan, et al., *Methods in Enzymol.*, 204:63-113 (1991)).

5           **E)     Assay for presence of desired cassettes**

10           The cassettes obtained will encode many different sorts of genes. In many cases, genes encoding functions of one particular kind but with differing specificities have related polypeptide sequences. A particular example of this kind of relationship is the set of genes that encode DNA methyltransferases, which carry out the same reaction (adding a methyl group to a specific base in a specific sequence) but with differing specificities (different particular bases within different particular sequences are modified). These can be tentatively identified by PCR employing primers that anneal to conserved polypeptide motif (Morgan, *supra* 15 (1996)). Briefly, individual colonies or pools of colonies from step D) can be subjected to degenerate PCR by procedures detailed in Morgan, 1996, with modification. Most suitable would be a design in which degenerate primers annealing to the methyltransferase motifs form one end of the amplicon and the other end of the amplicon is formed by one or more of the primers annealing to the flanking repeats. If a PCR product of suitable size is obtained, the relevant colony is likely to contain a gene for a methyltransferase. Plasmid or phage clones from candidate colonies identified in this way can then be sequenced in part or in whole.

25           Alternatively, plasmid or phage clones from colonies picked at random can be sequenced. Clones with potential methyltransferase genes can be identified by evaluation using DNA comparison algorithms such as BLAST or FASTA, or by means of programs specifically directed to evaluating such similarities (Posfai, et al., *Compt. Appl. Biosci.* 10(5):537-544 (1994)).

30           Functional tests for specific activities can also be use, as in Example 2.

## EXAMPLE 2

### FINDING RESTRICTION ENZYME CASSETTES BY FUNCTIONAL REPORT FOLLOWED BY CHARACTERIZATION

5 The present procedure will allow isolation in expression-ready form of a large number of cassettes specifying a variety of genes with diversity-selected functions. Accordingly, identification of specific clones expressing functions of the desired type is a critical part of the procedure. This example illustrates one way  
10 to identify a particular desired function, a DNA damaging agent, and to refine the functional identification until a site-specific doublestranded DNA endonuclease (a restriction enzyme) has been characterized. In addition, this example illustrates that the method is useful even when the desired function is toxic to the cell that expresses it. The procedure of this Example is possible specifically because the  
15 orientation of the genes is specified in advance, due to the natural orientation of the genes in a cassette array relative to the repeat elements that separate them.

20 Accordingly, in one embodiment, the vector employed, pLT7K (Fig 10), can be used to regulate the expression of the cloned cassette fragments even when nothing whatever is known about the identity or sequence of the cassettes individually. In this vector, two levels of control are available: expression is inducible and inhibition is repressible. A T7 gene 10 promoter reads into one side of the cloning site; expression from this promoter is repressed by LacI provided by the vector, as is expression of the T7 RNA polymerase itself, which is provided by  
25 the host cells used for expression. Further control can be obtained by the use of pLysP, which expresses an inhibitor of T7 RNA polymerase.

30 To further reduce expression directed by the cloned fragment, and residual leaky expression from the T7 promoter, tandem  $\lambda$  pL promoter reads into the other side of the cloning site, antagonizing expression from pT7. This antagonistic transcription is regulated by  $\lambda$  cI<sup>857</sup>, a thermosensitive repressor. At 40°C and in the absence of IPTG therefore, essentially no expression was observed; at 30°C, some leaky expression is seen; at 30°C in the presence of IPTG, moderate levels of expression can be achieved.

The strategy employed in the present Example, an indirect report of DNA damage is used to identify those cloned cassettes that lead to DNA damage, a procedure carried out by subjecting a portion of each clone to conditions that induce expression of the cassettes, and examining the color of colonies thus induced. Those that yield a positive signal are then chosen, and the portion of the clone never subjected to the inducing condition is carried to the next step. This ensures that the DNA damage step does not select for inactivation of the gene identified. The positive cassettes identified at this step (a reduced number) can then be examined in more detail. These are then examined by inducing another portion of each clone and examining the induced portion for three indices of site-specific DNA cleavage. Finally, the clones of interest are sequenced.

#### A. Reporters of DNA damage for use with pLT7LK.

In order to use the DNA damage indicator strategy for identification of DNA damaging cassettes cloned into pLT7LK, a host strain was required with five characteristics: the T7 RNA polymerase should be expressible after induction; the strain should not contain a lambda lysogen (because it would be induced to express phage-encoded killing functions following DNA damage); it should preferably be highly transformable, in order to obtain a large collection of transformants carrying cloned cassettes; it should express the DNA damage indicator *lacZ*, preferably only following DNA damage--ie with a clean background of white colonies in the absence of induction; and it should not express the major nonspecific endonuclease of *Escherichia coli*, Endonuclease I. This last requirement is needed for clear identification of restriction digest banding patterns in agarose gels, resulting from the action of site-specific endonucleases on test DNA substrates.

ER2745 and ER2746 were constructed by standard P1vir transduction. These strains provide alternative host backgrounds with differing advantages, both useful for the present goal of identifying cassette clones in pLT7K that cause damage to DNA when expressed.

A sample of the ER2745: (F  $\lambda$  *fhuA2* [*lon*] [*dcm*] *ompT lacZ::T7 gene1* *gal sulA11*  $\Delta$ (*mcrC-mrr*)114::*IS10 R(mcr-73::miniTn10--TetS)2 R(zgb-210::Tn10* --TetS) *endA1*) *dinD2::MudI1734* (*KanR, lacZ<sup>+</sup>*) has been deposited with the American Type Culture Collection under the terms and conditions of the Budapest Treaty on \_\_\_\_\_, 1999 and has received ATCC Patent Deposit No. \_\_\_\_\_.

A sample of ER2746: (F  $\lambda$  *fhuA2 glnV44 e14- rfbD1? relA1? endA1* *spoT1? thi-1*  $\Delta$ (*mcrC-mrr*)114::*IS10 lacZ::T7 gene1 dinD2::MudI1734* (*KanR, lacZ(ts)*) has been deposited with the American Type Culture Collection under the terms and conditions of the Budapest Treaty on \_\_\_\_\_, 1999 and has received ATCC Patent Deposit No. \_\_\_\_\_.

ER2745 was constructed in one step from an existing strain. The existing strain, ER2566, was deficient in all known endogenous restriction systems (enabling efficient cloning), did not express  $\beta$ -galactosidase, and expressed T7 RNA polymerase under *lacI* control from a chromosomal location (not an inducible prophage). It also lacked Endonuclease I, the major nonspecific nuclease of *E. coli*, and so would be useful for visualizing restriction enzyme activities in crude extracts. The *dinD* indicator was introduced into this strain by P1 transduction from strain ER1992 of Fomenkov, *supra* (1995)), to form ER2745.

ER2746 was constructed in three steps from an existing strain. The existing strain, ER2418, had the desirable property of relatively high induced competence, a property shared by many lines derived from *E. coli* K12 but not present in lines derived from *E. coli* B like ER2745. The allele for expression of T7 RNA polymerase was introduced in two transductional steps: ER2418 x P1(ER2556) --> Tet<sup>R</sup> (Pro- Kan<sup>R</sup>) to form ER2740; then ER2740 x P1(ER2553) --> Pro<sup>+</sup> (Kan<sup>S</sup> Tet<sup>S</sup> Lac- T7RNAP<sup>+</sup>) to form ER2744. Finally, a *dinD* indicator allele was introduced into ER2744 from ER2170.

## B. Cloning the cassettes

Cloning of cassettes was carried out by amplification from chromosomal samples. Total genomic DNA of *P. alcaligenes* (ATCC No. 55044) (NEB



Deposit No. 585, New England Biolabs, Inc., Beverly, MA) prepared by the procedure of Qiagen (Genomic tip 100/G (Cat 10243) as above was amplified using 8 combinations of primers 8-13 (SEQ ID NO:86 through SEQ ID NO:91 respectively; see Table 2): 8+12, 9+12, 10+12, 11+12 and 8+13, 9+13, 10+13, 11+13. The various combinations enable efficient amplification from different families of PAR repeat elements, since the central portion within each family of oligonucleotides (8-11 or 12-13) is varied in sequence. Each of the different versions facilitates annealing to different family members.

PCR amplification was by the procedure of Example 1, Section C2. Amplified cassettes were then digested with 20 units XbaI and 1 unit XhoI (New England Biolabs Cat. Nos. 145 and 146, Beverly, MA) in 1X NEBuffer 2 for 1 h at 37°C. Digested fragments were ligated overnight at 16°C with doubly-digested, dephosphorylated pLT7K. Dephosphorylation was for 1 h at 37°C with shrimp alkaline phosphatase (Amersham #E70092Y); ligation was with NEB Catalog No. 202 (New England Biolabs, Inc., Beverly, MA). These ligated libraries were introduced into ER2745 and ER2746 by electroporation, followed by plating on LB + ampicillin (100 µg/ml) and incubation overnight at 40°C. At this temperature, antisense expression is derepressed and in the absence of IPTG sense expression is uninduced, yielding expression undetectable by the DNA damage indicator described below (Section C).

### C. Screening for functional report.

The clone library thus recovered under conditions that repress expression of the integron cassettes (40°C -IPTG) to assure viability can then be scored for functional report. Replica plating onto Xgal plates and incubation under semi-inducing (30°C) or inducing (30°C +IPTG) conditions will allow identification of colonies that express DNA damaging functions. Some of these will be restriction enzymes. Individual colonies can then be recovered from master plates that have not been subjected to the damaging condition, to assure recovery of the original sequence.

## D. Assessment of clone identity

The DNA damage screen can allow identification of RM genes (Fomenkov, *supra* (1995); Fomenkov, *supra* (1994)). However, other sorts of genes will also be obtained; for example, a single-strand specific nuclease was among the genes recovered using the Endo-Blue method (Fomenkov, *supra* (1994)). Three procedures can be used to identify RM genes. In the first, cells are induced to express the cassette-encoded genes, crude extracts are made, these extracts are used to digest standard target DNAs, and enzymatic activity is detected by production of discrete bands on agarose gels. In the second, clones are briefly induced to express the cassette-encoded gene, then the whole cells are subjected to pulse-field gel analysis. Discrete bands will result from digestion of the chromosomal DNA of the clone-bearing cells. In the third approach, sequencing of clones to allow classification by homology searches.

### D1) Crude extract assay

Clones positive in the DNA-damage screen will be grown under non-inducing conditions to late log phase, and shifted to the inducing condition for four hours. This procedure was successful in allowing expression of an amount of PacI similar to that expressed in the native host, *P. alcaligenes* (D. Byrd, personal communication). Cells are collected by centrifugation, resuspended in buffer, lysed by lysozyme-EDTA treatment, clarified by centrifugation.

Digests are of three sorts:

- 1) a PacI-specific digest using a specific substrate designed to give a diagnostic pattern, for the positive control.
- 2) a general screen for 4-6 base cutters, using standard plasmid, phage and viral DNAs. Some 8-base specificities may be detected by this method as well.
- 3) a general screen for 8-base cutters. In vitro screens for enzymes with 8-base sites are more difficult because of the rarity of sites. However, it is usually possible to distinguish between nonspecific nuclease and an 8-base endonuclease

using total chromosomal DNA as a substrate for in vitro digestion with crude. This is due to the presence of specific fragments (especially large ones) not subject to further digestion; even though the fragments are not resolvable on the gel (and the recognition site cannot be deduced), the result is recognizably different from that produced by nonspecific nucleases (which preferentially degrade large fragments).

In each case, aliquots of extract are incubated with potential DNA substrates in the presence of  $Mg^{++}$ . Products will then be analysed by agarose gel electrophoresis.

#### D2) Pulsed-field gel assay

A potentially more-informative assay for 8-base recognition sites would rely on separation of total chromosomal fragments on pulsed-field gels. When crude extracts are used for screening procedures, these gels are too cumbersome and too sensitive to other nucleases in the extract to be generally useful. However, in this case we can adapt the procedure to our purposes

In standard procedures, the substrate DNA is obtained by first embedding whole cells in agarose plugs. DNA is released from the cells in situ by means of a series of enzymatic treatments and washes that degrade the cell wall. The restriction endonuclease is then incubated with the plug; this usually takes several hours, since the enzyme must permeate the agarose and the remnants of the previous digestions.

The restriction nuclease digestion step can be bypassed by inducing expression within the cell, before agarose is added. By definition, the candidate clones are known to damage DNA in vivo in regulated manner. Accordingly, a banding pattern should be identifiable using the chromosomal DNA of the cells in which expression of the enzyme is induced. *PacI* will again be used as a test case. *NotI* will also be used, since the pattern expected for a total chromosomal digest is already well-known.

Critical steps are: quenching endogenous DNA degradation (especially exonuclease activity) at harvest and during the agarose-embedding process; the

length of the induction; and the degree of induction. Controls include: positive control, standard digestion of the host DNA embedded in agarose plugs with purified PacI and NotI; and negative control, samples of the host containing the empty vector, treated in parallel with the experimental samples.

Improvements in the strain used for this part of the survey include introduction of a *recD* mutation, which would inactivate the major ATP-dependent double-strand exonuclease of the cell; and introduction of an *xth* mutation that would inactivate the major ATP-independent double-strand exonuclease. A triply nuclease-deficient strain (*endA xth recD*) should be viable but may not stably maintain the plasmid (Niki, et al., *supra* (1990)).

### D3) Sequencing

Genes obtained can be sequenced. To reduce redundant sequencing efforts, restriction digestion and fingerprinting of large numbers of candidates can be carried out. The recovered genes into sets with similar fingerprints, and two of each are sequenced. A minimum of three-fold sequence coverage is usually required in order to have sufficient confidence to carry out preliminary homology searches.

Sequencing can be conducted efficiently using the newly available Tn7-based transposition system, GPS™-1 (New England Biolabs Catalog No. 1700, New England Biolabs, Inc., Beverly, MA). This system enables introduction of primer-binding sites at random locations in plasmids of interest, rapid mapping of the location of the insertion by digestion with rare-cutters that cleave within the transposon, and sequencing of the insertions within the fragment of interest. With these target molecules, About 20% of transposon insertions will be found within the sequence of interest. No more than 6 suitable insertions are needed in most cases, since cassettes are normally smaller than 2 kb. Two sequence runs (500 bp per run) from flanking vector primers and 12 runs from insertions will yield 7000 bp of raw sequence, approximately 3-fold redundancy. This is be sufficient for primary analysis. Further sequencing can be carried out to obtain high-quality sequence of the most interesting fragments. Other sequencing strategies are also possible.

Homology to genes in public databases can help to exclude candidates for new type II RM genes. Many genes that might be recovered during this procedure exhibit conserved amino acid sequence segments: topoisomerases, helicases, nicking enzymes associated with conjugal plasmid transfer, and transposases all can be found annotated in databases, identified by BLAST or other homology search procedures. Genes for type II restriction enzymes, on the other hand, rarely can be identified in this way. When they can be identified by homology, they are almost always isoschizomers of (recognize the same site as) the enzyme in the database (R. Roberts, personal communication). Thus, the target genes (endonucleases recognizing new specificities) can be expected among those not identified by homology search.

These target genes, for type II endonucleases of unknown specificity, normally can best be identified by adjacency to genes encoding protective modification methyltransferases (R. Roberts and J. Posfai, personal communication). Methyltransferases are recognizable by bioinformatic methods, since conserved motif elements are always present (see above). However, two enzymes that should be recoverable by the present method, *PacI* and *PmeI*, are not adjacent to genes similar to any modification methyltransferase, and indeed so far no protective methyltransferases have been identified in the original hosts. Since these enzymes recognize AT-rich 8-base sites and the host organisms contain GC-rich genomes, host protection may be achieved by means of absence of sites.

Accordingly, candidate type II endonuclease genes of special interest will be solo ORFS with no database hits. Candidates adjacent to identifiable methyltransferase genes will be also retained, as will potential isoschizomers, which could have other desirable properties such as those affecting stability.

### EXAMPLE 3

#### GENERAL PROCEDURE FOR EMPLOYMENT OF THE METHOD

Repeats to be sought include those in the public literature (Hall and Stokes, *Genetica* 90:115-132 (1993); Hall and Collis, *Mol Microbiol* 15:593-600 (1995);

Levesque, et al., *Gene* 142:49-54 (1994); Recchia and Hall, *Mol Microbiol* 15:179-187 (1995); Mazel, et al., *Science* 280:605-608 (1998); Barker, et al., *J Bacteriol* 176:5450-5458 (1994); Clark, et al., *Mol Microbiol* 26:, 1137-1138 (1997); Ogawa and Takeda, *Microbiol Immunol* 37:607-616 (1993); Hall, et al. *Mol Microbiol* 5:1941-1959 (1991); Levesque, et al., *Antimicrob Agents Chemother* 39:185-191 (1995); Sallen, et al., *Microb Drug Resist* 1:195-202 (1995); Sandvang, et al., *FEMS Microbiol Lett* 160:37-41 (1998); Senda, et al., *J Clin Microbiol* 34:2909-2913 (1996); Tosini, et al., *Antimicrob Agents Chemother* 42:3053-3058 (1998)) those disclosed herein (SEQ ID NO:5 through SEQ ID NO:74), and those identified in the genome sequence of one or more model organism of interest. The set of repeat sequences identified in the organism of interest are determined by the method of Example 1. These segments are then made into a multiple alignment, for example using the program MEGALIGN (DNASTAR, Madison Wisconsin) and preferably the CLUSTAL method of alignment within it. Segments thus identified can be grouped into families, for example by means of the Phylogeny facility in the MEGALIGN program, and bushy groups, in which there are many interior branches, are chosen as repeat families. These additional families should direct the design of oligonucleotides for use as probes or primers during application of the method.

## 2) Identification of a variable class of functions

A function of interest is identified in a taxon related to the model organism of interest. This can be for example ability to adhere to a particular tissue, for example red blood cells or the root hairs of plants.

A relatively large (>50 members) and diverse collection of isolates within the taxon of interest are collected. The diversity of these isolates is characterized by isolation from locations spanning the extremes of the organism's distribution; these extremes may include spatial (geographic) distribution, thermal tolerance, salt tolerance, pH tolerance, O<sub>2</sub> partial pressure tolerance or requirement or host organism identity.

The members of this collection are screened for the presence of the function of interest and its specificity. In this example, it may be done by testing for

hemagglutination ability, with red blood cells of sheep, cows, rabbits, pigs, goats, frogs, and humans as examples of different specific targets, or may be tested with one type of red cell in the presence of different mono- or disaccharides, or following various treatments that alter the nature of the red cell surface. The function is identified as variable in the way that is expected of cassette-encoded functions if one or both of two conditions obtains. First, a large fraction ( $>10\%$ ) is different from the rest, in whether the function is present or absent. For example, 5 or more members of the collection express hemagglutination of the red cells, and the rest don't; or vice versa. Second, the specificity of the function varies: for example, some agglutinate sheep red cells, others goat red cells. This criterion is best satisfied if the number of specificities identified is large, for example  $>4$  different specificities in a collection of 50 isolates.

Variable functions can also be identified by immunological procedures, for example ELISA assays employing sera from animal or human populations of interest, or monoclonal antibodies recognizing variable epitopes in a compound of interest (e.g. a polypeptide); or by cytotoxicity assays, for example employing tissues of different physical or phylogenetic origins; or assays testing inhibition or stimulation of cellular processes such as DNA synthesis or cAMP hydrolysis directly or indirectly, in a context of tissue- or organism-specific effects; or tests of growth on or transformation of varied potential sources of carbon, nitrogen, or energy; or tests of growth in the presence of or inhibition of varied antimicrobial compounds.

### **3) DNA preparation and determination of suitability for use of the method**

A preliminary test of the suitability of the method may be carried out by colony PCR, by inoculating a series of small samples of culture medium (for example in microtiter well plates) with portions of isolates of the taxon to be examined (reserving another portion for storage), growing them, boiling them, and carrying out PCR as in Example 1, Part C2. Other primers designed based on these or other repeat families identified from the literature or in step 1 can also be used. Positive isolates identified at this step by the appearance of one or more PCR product are then carried to the next step.

#### 4) Cassette isolation

DNA preparations from positive isolates is subjected to PCR on a larger scale, employing primer pairs with suitable restriction enzyme cloning sites at the ends as in Example 2: SEQ ID NO:86 with SEQ ID NO:90; SEQ ID NO:86 with SEQ ID NO:91; SEQ ID NO:87 with SEQ ID NO:90; SEQ ID NO:87 with SEQ ID NO:91; SEQ ID NO:88 with SEQ ID NO:90; SEQ ID NO:88 with SEQ ID NO:91; SEQ ID NO:89 with SEQ ID NO:90; SEQ ID NO:89 with SEQ ID NO:91 (see Table 2). Additional primer pairs designed based on additional repeat families may also be designed. Amplification conditions may be adjusted depending on the pairs used.

#### 5) Cassette cloning

The PCR fragments are digested with XhoI and XbaI if the primers of Example 2 and pLT7K are used; other primers can be used including primers suitable for use with a derivative of pLT7K or similar plasmid carrying other restriction sites at the cloning site.

#### 6) Strain choice

A strain suitable for recovery of cassettes will be one not expressing the function of interest, but in which its presence can be sought. For example, hemagglutinin genes should be expressed in a strain not itself expressing a hemagglutinin that would interfere with the survey. LE392 is an example of an E. coli strain that does not express hemagglutinin activity. For use with pLT7K, the T7 gene1 construct would need to be introduced into LE392; or alternatively, strains such as ER2645, ER2746, ER2566 or ER2744 could be used if they were shown to lack hemagglutinin activity. The strain may be customized to facilitate expression or report of functionality, for example by expressing a protein export system capable of exporting a class of hemagglutinins sought (eg. fimbriae).



## 7) Cassette identification

In the case of hemagglutination, a functional assay is available, so colonies or pools of colonies can be tested for hemagglutination in microtiter wells, following induction of expression as in Example 2.

Another method of identification would be to design degenerate primers specific for motifs found in particular classes of expected proteins, for example fimbriae, pili, or outer membrane proteins, and use them to perform PCR on colonies or pools of colonies either alone or in combination with PCR primers specific for the flanking repeats, as described in example 2.

A list of motifs characteristic of classes of proteins can be found in the public databases described in (M. Patterson and M. Handel, "Trends Guide to Bioinformatics" Elsevier Science, Cambridge, UK, (1998)).

## 8) Functional characterization

Colonies specifically exhibiting properties expected of desired gene cassettes would then be characterized by methods appropriate to the particular function identified, for example, in a hemagglutination test by competition with small molecules such as various sugars; by its sensitivity to various treatments such as iodination, heating, freezing, treating with acid, alkali, or alkylating agents or with proteases or nucleases; and by obtaining the sequences of the genes and determining the properties of cells with genes carrying mutations of various sorts including fusions to other reporter molecules such as alkaline phosphatase, beta galactosidase, green fluorescent protein or various epitope tags, or obtaining purified preparations of encoded proteins by standard purification methods or by affinity purification by means of polypeptide tags.

**WHAT IS CLAIMED IS:**

1. A method for the cloning of intact, diversity-selected genes from within  
gene cassettes, said method comprising the steps of:

(a) identifying repeat DNA sequences which flank gene cassettes;

(b) hybridizing oligonucleotides to said repeated sequences which flank  
said gene cassettes and amplifying said sequences to provide DNA fragments  
which contain genes from within the cassettes.

(c) ligating said DNA fragments into a vector; and

(d) transforming said vector into an appropriate strain.

2. The method of claim 1 wherein said diversity-selected genes are selected  
from the group consisting of:

cell surface antigens such as polysaccharide antigens or polypeptide  
antigens or secreted molecules; adhesins such as fimbrial proteins, pilus proteins or  
outer membrane proteins; transporters of small molecules, especially those with  
narrow specificity; toxins, hemolysins, hemagglutinins, kinases and signaling  
molecules;

detoxifying enzymes such as drug resistance determinants; catabolic  
enzymes specific for compounds episodically available, excluding those required  
for central metabolic pathways such as the tricarboxylic acid cycle; enzymes for  
biosynthesis of rare sugars, excluding those required in all cells, such as ribose,  
deoxyribose, and sugars of the cell wall, especially of those sugars that form part  
of the pericellular envelope.

3. The method of claim 2 wherein said diversity-selected genes comprise  
restriction endonuclease genes.

4. The method of claim 2 wherein said diversity-selected genes comprise  
methyltransferase genes.

5. The method of claim 1 wherein said oligonucleotides contain recognition sites which permit directional cloning.

5 6. The method of claim 5 wherein the DNA fragments are ligated into said vector in an orientation that enables expression.

7. A method for identifying the presence of gene cassette arrays from within a target DNA preparation, said method comprising the steps of:

10 (a) hybridizing at least one oligonucleotide which hybridizes to one or more of SEQ ID NO:5 through SEQ ID NO:78 to a DNA preparation; and

(b) detecting the presence of a stable DNA-DNA hybrid.

15 8. The method of claim 7 wherein said detection comprises determining the presence of stable DNA-DNA hybrid by Southern blot or dot blot.

20 9. The method of claim 7 wherein said detection comprises employing at least two oligonucleotides and hybridizing said oligonucleotides to said DNA preparation, and detecting their ability to support DNA polymerization at the 3' end of the stable DNA-DNA hybrid.

25 10. The method of claim 7 wherein said oligonucleotides comprise SEQ ID NO:79 through SEQ ID NO:91.

11. The method of claim 7 wherein said oligonucleotides hybridize to one or more of DNA SEQ ID NO:5 through SEQ ID NO:78 or portions thereof.

30 12. The method of claim 7 wherein the DNA source comprises an individual strain.

13. The method of claim 7 wherein the DNA source comprises a group or pool of strains.

14. The method of claim 7 wherein the DNA source comprises environmental DNA.

5 15. A composition consisting of isolated DNA primers comprising SEQ ID NO:79 through SEQ ID NO:91 or portions thereof.

16. A composition consisting of DNA primers which hybridize to one or more of DNA SEQ ID NO:5 through SEQ ID NO:78 or portions thereof.

10 17. A method for identifying gene cassette arrays from a predetermined DNA sequence, said method comprising the steps of:

(a) screening the said predetermined DNA sequence for TAACWA;

15 (b) screening the said predetermined DNA sequence for CGTTRR;

(c) screening for DNA segments wherein the 5' T of step A is less than about 200 base pairs from the 3' R of step B; and

20 (d) determining whether the DNA sequence of step C is repeated in the predetermined DNA sequence.

09/701 626

Table 1

SEQ ID NO:5 through SEQ ID NO:39

SEQ ID NO: #	Sequence
SEQ ID NO:5	ATCTAACAAAT TGGTTCAAGT CGCTCGCTTC GCTCACTCGG GACCGGCTAA AGCCGGCCCC TTAAACCAAGC GTTAGGT
SEQ ID NO:6	ACCTAACATG GCGTCAACC GGGCTCCCTT CCGTCGCTGG AGCTGGCGG ATAAAGCCGC GCAGCGCGG TTAGCTCTAC GTTAGGG
SEQ ID NO:7	GTCTAACAAAT TGGTCAAGT CGTTCGCTTC GCTCACTCGG GACTCCGCA AGCTCGCTC GCGCGCGCC CTAGCCCAA CGTTAGG
SEQ ID NO:8	ACCTAACAAAT GCGTCAACT GCGCTCACT CCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:9	ACCTAACAAAT GCGTCAACT GCGTCACT TCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:10	GCCTAACAAAT GCGTCAAG CCGTCACTTC GTTCGCTGG ACCGCGAAG CCGGCCCTT AGCTTAATCG TTAGGT
SEQ ID NO:11	ACCTAACAAAT TGGTTCAAGT CGTTCGCTTC GCTCACTCGG GACCGGCTAA AGCCGGCCCC TTAAACCAAC GTTAGGC
SEQ ID NO:12	GCCTAACAAAT TGGTCAAGT CGTTCGCTTC GCTCACTCGG GACCGGCTAA AGCCGGCCCC TAGCCAAAGC TTAGGT
SEQ ID NO:13	ACCTAACAAAT GCGTCAACT GCGTCACT TCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:14	GCCTAACAAAT TGGTCAAGT CGTTCGCTTC GCTCACTGG ACGGCTAA GCGGCCCTT TAACCAACG TTAGGC
SEQ ID NO:15	CTCTAACAAAT GCGTCAACT ATCGTCACT TCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:16	TTATAACAAT GCGTCAACT GCGTCACT GCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:17	ATTTAACAAAT GCGTCAACT GCGTCACT GCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:18	CCCTAACAAAT TGGTCAAGT GCGTCACT GCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:19	CTCTAACAAAT TGGTCAAGT GCGTCACT GCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:20	GCCTAACAAAT TGGTCAAGT GCGTCACT GCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:21	GCCTAACAAAT TGGTCAAGT GCGTCACT GCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:22	CCCTAACAAAT TGGTCAAGT GCGTCACT GCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:23	GCCTAACAAAT TGGTCAAGT GCGTCACT GCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:24	ACCTAACAAAT TGGTCAAGT GCGTCACT GCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:25	GCCTAACAAAT TGGTCAAGT GCGTCACT GCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:26	GCCTAACAAAT TGGTCAAGT GCGTCACT GCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:27	ACCTAACAAAT TGGTCAAGT GCGTCACT GCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:28	GCCTAACAAAT TGGTCAAGT GCGTCACT GCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:29	ACCTAACAAAT TGGTCAAGT GCGTCACT GCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:30	GCCTAACAAAT TGGTCAAGT GCGTCACT GCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:31	GCCTAACAAAT TGGTCAAGT GCGTCACT GCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:32	GCCTAACAAAT TGGTCAAGT GCGTCACT GCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:33	GCCTAACAAAT TGGTCAAGT GCGTCACT GCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:34	GCCTAACAAAT TGGTCAAGT GCGTCACT GCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:35	GCCTAACAAAT TGGTCAAGT GCGTCACT GCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:36	GCCTAACAAAT TGGTCAAGT GCGTCACT GCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:37	GCCTAACAAAT TGGTCAAGT GCGTCACT GCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:38	GCCTAACAAAT TGGTCAAGT GCGTCACT GCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG
SEQ ID NO:39	GCCTAACAAAT TGGTCAAGT GCGTCACT GCGTTCGCTG GACTCGCAA AGCTCGCTT TTGCTCGCC GTTAGCTTAA TCGTTAG

Table 1 Continued  
SEQ ID NO:40 through SEQ ID NO:75

SEQ ID NO:40	GCCTAACAAAT AGGTTCAAGT GCTCGCTTC GCTCACTGG GACCGGCTAA AGCCGGCCCC TTAACCAAAC GTTAGGT
SEQ ID NO:41	ACCTAACTAG TGGTTCAAGC CGTCTGCTTC GCTCACTGG GACCGGCTAA AGCCGGCCCC TTAACCAAAC GTTAGGC
SEQ ID NO:42	ATATAACAAT TGGTTCAAC CGTTGCTGC GCTCACTGG ACGGCTAA GCGGCCCT TAACCAAACG TTATGC
SEQ ID NO:43	GCATAACAAT TGGCTAAGC CGTCTGCTCC GCTCACTGG AGTCTGCTAA GCTACGTTTC CGCGGCCCTC TTAGCCAAAC GTTAGGC
SEQ ID NO:44	CCCTAACAAA TGGTTCAAG CGTTGCTTC CGTCACTCG GGACCGGCTA AGCCGGCCCC CTTAACCAAAC GTTAGGC
SEQ ID NO:45	GCCTAACAAAT CGGCTCACT CGGCTCACT TCGTCTGCTG GACAGTCAA AGCTGCTCT TTGCTGCTC GTTAGCTTAA TCGTTAG
SEQ ID NO:46	CTCTAACAAAT GCGCTCAAT CGTCACTAC GTTCGCTGG ACCGCTAA GCGGCCCTC TAGCTTATC GTTAGAG
SEQ ID NO:47	CTCTAACAAAT TGGTTCAAGT CGTTGCTTC GCTCACTGG GGACCGGCTA AGCCGGCCCC CTTAACCAAAC GTTAGGG
SEQ ID NO:48	CCCTAACAAA TGGTTCAAGT CACTGCTTC GCTCACTGG GACCGGCTA AGCCGGCCCC TTAACCAAAC GTTAGAG
SEQ ID NO:49	CTCTAACAAAT TGGTTCAAGT CGTTGCTTC GCTCACTGG GGACCGGCTA AGCCGGCCCC CTTAACCAAAC GTTAGGC
SEQ ID NO:50	ATCTAACAAAT TGGCTCAAGT CGTTGCTTC GCTCACTGG GACGCTCAAT AGTCTGCTA TTGGCGGCC CTTAGCCAAAC GTTAGGC
SEQ ID NO:51	GCCTAACAAAT TGGTTCAAGT CGTTGCTTC GCTCACTGG GGACCGGCTA AGCCGGGCC CTTAACCAAAC GTTAGGG
SEQ ID NO:52	CCCTAACAAA TGGTTCAAGT CGTTGCTTC GCTCACTGG GACCGGCTA AGCCGGGCC CTTAACCAAAC GTTAGGC
SEQ ID NO:53	GCCTAACAAAT TGGTTCAAGT CGTTGCTTC GCTCACTGG GGACCGGCTA AGCCGGGCC CTTAACCAAAC GTTAGGC
SEQ ID NO:54	GCCTAACAAAT TGGTTCAAGT CGTTGCTTC GCTCACTGG GACGCTCAAT AGTCTGCTA TTGGCGGCC CTTAGCCAAAC GTTAGGC
SEQ ID NO:55	TTCTAACAAAT TGGTTCAAGT CGTTGCTTC GCTCACTGG GACCGGCTA AGCCGGGCC CTTAACCAAAC GTTAGGC
SEQ ID NO:56	GGCTAACAAAT CGGCTCACT GTGCTCACT TCGTCTGCTG GACAGTCAA AGCTACGCTT TTGCTGCTC GTTAGCTTAA TCGTTAG
SEQ ID NO:57	GCCTAACAAAT TGGTTCAAGT CACTCACTTC GCTCACTGG GACCGGCTC GCGGCCCTC TAACCAAACG TTGGGC
SEQ ID NO:58	GCCCAACAAT CGGCTCACT GCGCTCACT TCGTTGCTG GACGCTCAA AGCTACGCTT TTGGCGGCC CTTAGCTTAA TCGTTAG
SEQ ID NO:59	GCATAACAAT TGATTCAGT CGTTGCTTC GCTCACTGG GGACCGGCTA AGCCGGGCC CTTAACCAAAC GTTAGGC
SEQ ID NO:60	GCCTAACAAAT TGGTTCAAGT CACTGCTTC GCTCACTGG GACCGGCTC GCGGCCCTC TAACCAAACG TTAGGC
SEQ ID NO:61	ATCTAACAAAT TGGCTCAAGT CGTTGCTTC GCTCACTGG GACCGGCTA GCGGCCCTC TAACCAAACG TTATGC
SEQ ID NO:62	GCATAACAAT TGGCTCAAGT CGTTGCTTC GCTCACTGG GACCGGCTA GCTACGCTT CCGCGGCCCTC TTAGCCAAAC GTTAGGC
SEQ ID NO:63	GCCTAACAAA TGGTTCAAGT CGTTGCTTC GCTCACTGG GACCGGCTA AGCCGGGCC CTTAACCAAAC GTTAGGC
SEQ ID NO:64	GCCTAACAAAT TCAGTCAAGT GCGGCTCAAC CCGGCTGCTC GGTCTTGG CCGCTTATC CAAGCGTTAG AT
SEQ ID NO:65	ATCTAACAAAT TGGTTCAAGT CGTTGCTTC GCTCACTGG GACCGGCTA AGCCGGGCC CTTAACCAAAC GTTAGAG
SEQ ID NO:66	CTCTAACAAAT TGGTTCAAGT CGTTGCTTC GCTCACTGG GACCGGCTA GATCGGCTA GCGGCCCTC TAACCAAACG TTAGAG
SEQ ID NO:67	CTCTAACAAAT TGGTTCAAGT CGTTGCTTC GCTCACTGG GACCGGCTA GATCGGCTA GCGGCCCTC TAACCAAACG TTAGGC
SEQ ID NO:68	GCCTAACAAAT TGGTTCAAGT CGTTGCTTC GCTCACTGG GACCGGCTA AGCCGGGCC CTTAACCAAAC GTTAGGC
SEQ ID NO:69	GTCTAACAAAT GCGCTCAAGT CGTTGCTTC GCTCACTGG GACCGGCTA GCGGCCCTC TAGCTTATC GTTAGCA
SEQ ID NO:70	TGCTAACAAAT GCGCTCAAGT CGTTGCTTC GCTCACTGG GATAGTCAA AGCTGCTCT TTGCTGCTC GTTAGCTTAA TCGTTAG
SEQ ID NO:71	GCCTAACAAAT TGGTTCAAT GCTGCTCC GCTGCTGG ACCGCTAT CCGGCCCTA ACCAAGCGTT AGAT
SEQ ID NO:72	ATCTAACAAAT TGGTTAAAC CGTTGCTTC GCTCACTGG ACCGCTAA GCGGCCCTC TAACCAAACG TTAGGT
SEQ ID NO:73	GTCTAACAAAT TGGTTCAAGT CGTTGCTTC GCTCACTGG GACCGGCTA ATTCGGGCC TTAGCAAACG TTAAGT
SEQ ID NO:74	ATCTAACAAAT TGGTTCAAGT CGTTGCTTC GCTCACTGG GACCGGCTA AGCCGGGCC CTTAACCAAAC GTTAGGC
SEQ ID NO:75	GTATAACAAT TGGTTCAAGT CACTGCTTC GCTCACTGG GACCGGCTA AGCCGGGCC CTTAACCAAAC GTTAGAT

Table 1 Continued  
SEQ ID NO:76 through SEQ ID NO:78

SEQ ID NO:76	GCATAACAAG TCGTCAAAAT CGTCACTTC GTTCGCTGG GACGGGCTAA AGCCGGCCCC TTAGCTTATC GTTAGGC
SEQ ID NO:77	GCCTAACAAA TGGTTCAGC CGTTCGCTTC GCTCACTCGG GACCGGCTAA AGCCGGCCCC TTAACCAAAC GTTAGAC
SEQ ID NO:78	CTCTAACAAA TGGTTCAGC CGTTCGCTTC GCTCACTCGG GACCGGCTAA AGCCGGCCCC TTAACCAAAC GTTAGAG

	1992	1993	1994	1995	1996	1997	1998	1999	2000	2001	2002	2003	2004	2005	2006	2007	2008	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018	2019	2020	2021	2022	2023	2024	2025	2026	2027	2028	2029	2030	2031	2032	2033	2034	2035	2036	2037	2038	2039	2040	2041	2042	2043	2044	2045	2046	2047	2048	2049	2050	2051	2052	2053	2054	2055	2056	2057	2058	2059	2060	2061	2062	2063	2064	2065	2066	2067	2068	2069	2070	2071	2072	2073	2074	2075	2076	2077	2078	2079	2080	2081	2082	2083	2084	2085	2086	2087	2088	2089	2090	2091	2092	2093	2094	2095	2096	2097	2098	2099	2100	2101	2102	2103	2104	2105	2106	2107	2108	2109	2110	2111	2112	2113	2114	2115	2116	2117	2118	2119	2120	2121	2122	2123	2124	2125	2126	2127	2128	2129	2130	2131	2132	2133	2134	2135	2136	2137	2138	2139	2140	2141	2142	2143	2144	2145	2146	2147	2148	2149	2150	2151	2152	2153	2154	2155	2156	2157	2158	2159	2160	2161	2162	2163	2164	2165	2166	2167	2168	2169	2170	2171	2172	2173	2174	2175	2176	2177	2178	2179	2180	2181	2182	2183	2184	2185	2186	2187	2188	2189	2190	2191	2192	2193	2194	2195	2196	2197	2198	2199	2200	2201	2202	2203	2204	2205	2206	2207	2208	2209	2210	2211	2212	2213	2214	2215	2216	2217	2218	2219	2220	2221	2222	2223	2224	2225	2226	2227	2228	2229	2230	2231	2232	2233	2234	2235	2236	2237	2238	2239	2240	2241	2242	2243	2244	2245	2246	2247	2248	2249	2250	2251	2252	2253	2254	2255	2256	2257	2258	2259	2260	2261	2262	2263	2264	2265	2266	2267	2268	2269	2270	2271	2272	2273	2274	2275	2276	2277	2278	2279	2280	2281	2282	2283	2284	2285	2286	2287	2288	2289	2290	2291	2292	2293	2294	2295	2296	2297	2298	2299	2300	2301	2302	2303	2304	2305	2306	2307	2308	2309	2310	2311	2312	2313	2314	2315	2316	2317	2318	2319	2320	2321	2322	2323	2324	2325	2326	2327	2328	2329	2330	2331	2332	2333	2334	2335	2336	2337	2338	2339	2340	2341	2342	2343	2344	2345	2346	2347	2348	2349	2350	2351	2352	2353	2354	2355	2356	2357	2358	2359	2360	2361	2362	2363	2364	2365	2366	2367	2368	2369	2370	2371	2372	2373	2374	2375	2376	2377	2378	2379	2380	2381	2382	2383	2384	2385	2386	2387	2388	2389	2390	2391	2392	2393	2394	2395	2396	2397	2398	2399	2400	2401	2402	2403	2404	2405	2406	2407	2408	2409	2410	2411	2412	2413	2414	2415	2416	2417	2418	2419	2420	2421	2422	2423	2424	2425	2426	2427	2428	2429	2430	2431	2432	2433	2434	2435	2436	2437	2438	2439	2440	2441	2442	2443	2444	2
--	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	------	---

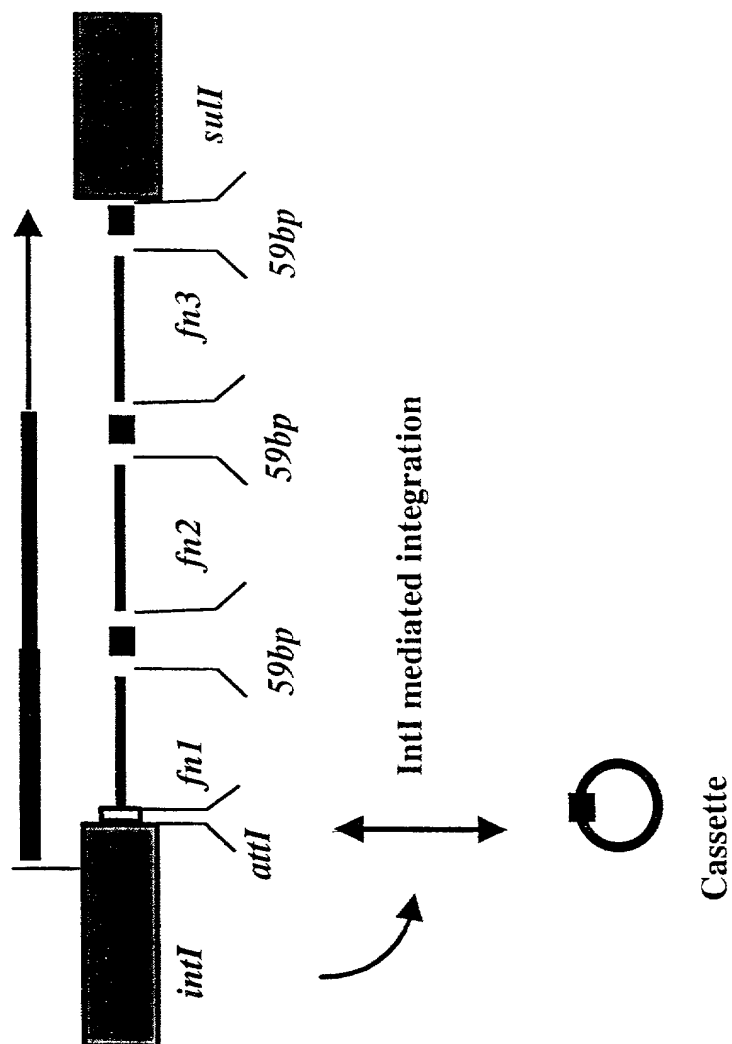
Table 2

SEQ ID NO:79 through SEQ ID NO:91

	Name in text	strand	Sequence
	SEQ ID NO:79	coding	TCGCTTCGGT CACTGGGGA CCGGCTAAAG CCGGCCCTT AACCAAAGT TA
	SEQ ID NO:80	coding	TAACAAATGG TTCAAGTCGT TCGCTTCGCT CACTGGGGA CCG
	SEQ ID NO:81	coding	TAACATATCA GTCAGCGGA CGCAACCCC GCTGCGGGT CTT
	SEQ ID NO:82	coding	TAACAAATGG CTCAACTGC CTCACTTCGT TCGCTGGACA GCC
	SEQ ID NO:83	coding	TAACAAGTCG CTCAACTGCC GCTCACTCGT TCGCTGGACA GCC
	SEQ ID NO:84	noncoding	GCCCCTTAAC CAACGTTA
	SEQ ID NO:85	coding	CCGAGTGAGC GAAGCAGCG
	SEQ ID NO:86	noncoding	AAACTCGAGG GTCCCCAGTG AGCGAAGCGA GCG
	SEQ ID NO:87	noncoding	AAACTCGAGG GTCCCGAGCG AGCGAAGCGA GCG
	SEQ ID NO:88	noncoding	AAACTCGAGG GTGTCCAGCG AGCGAAGCGA GCG
	SEQ ID NO:89	noncoding	AAACTCGAGA CCGCGCAGCG GGGTTTSCGT CCG
	SEQ ID NO:90	coding	TGCTCTAGAC GGCCCCCTTAA CCAACGTTA G
	SEQ ID NO:91	coding	TGCTCTAGAC GGCCGCCCGT TAGCTTAATC GTTAG

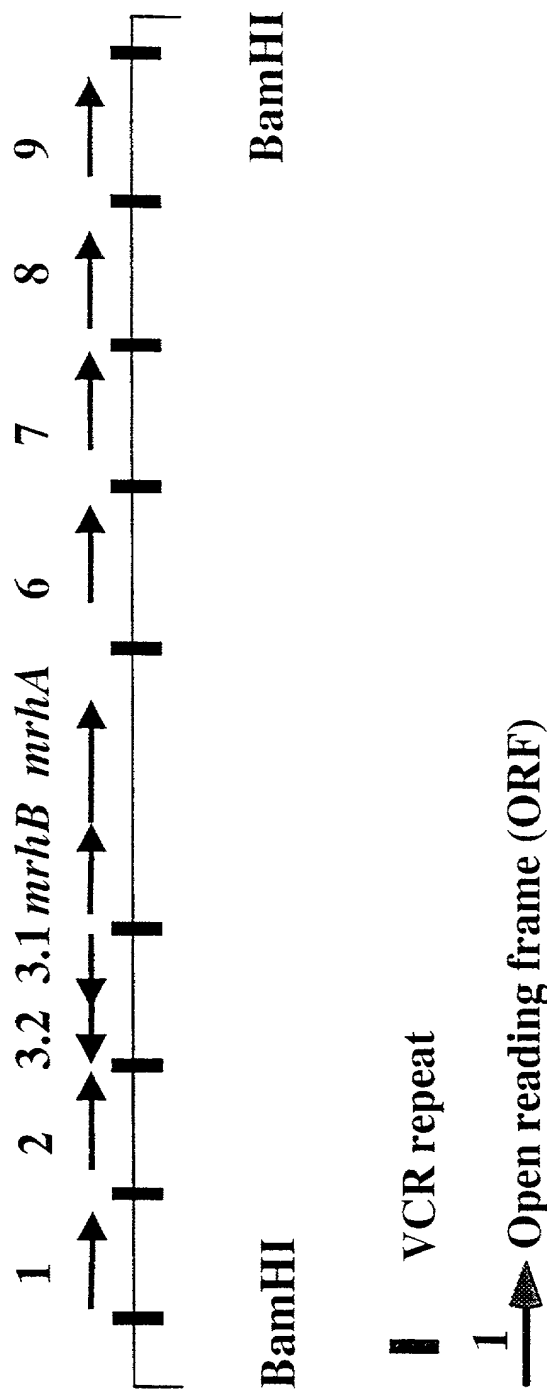


Fig 1 Integron structure



after Hall and Collis 1995

Fig. 2  
*Vibrio cholerae* superintegron  
fragment carried on pPM147



ORF2 is similar to Vlp()

ORF3.1 is similar to a plasmid-encoded protein (gi|516610)

ORF3.2 is similar to RelE (gi|42701) and plasmid-encoded proteins

[illegible]

CATACTTTCC AGGTTGTCTT TCTACGAAAA AATTGAAATC CTGAGAAAAA  
TTGATTTTGG TAACAGTCTC AAATCCCAAG AAAATACAGC GCTGCACCTA  
GACAAACTGA GCGGATTGCG TAACGCATTG GCGCATGCAG CACACATGCC  
ACCTGATGAA ATCATGAAGT TGTGCTCTGA TAAGTGGATA GAGTCCTTTG  
TGCTCGGATA TCCAAAGTCC ATTGGCAAAG AGAAAAATGC ACTTGAAAA  
CGGCTATCAC TTCTGTGGAA TTA CTGCCAC AGGAGGCATG TAGCAAAAA  
TAAGCAGCTT GCACACGAAC TCAAAAAATAC AGAGCAAGCC AACTAATAGA  
GTCCAGTTAT ACAGGTCCGT AAATGAGCCG CCTAACAACT GGTTCAGGCC  
ACTCACTTCG CTCGCTCGGG ACCGCGTTCC GCGGCCCTT AACCAAACGT  
TGGGCACCCA TAGAAAAATC CTAATGAGAA AACTATTCTAT ACCACTAATT  
TTCGCCCTGC TATCGGAGAG CTGATGGCA TCTGAAGCGT ATAAGGACCT  
TGAAACACAA GTAAGTAAA AAGCCAGCCT AGCAGTTGCC CAAATGAATG  
ACAGAGCAAC TGGAAAGCTC GACTACTCGG AAGAAAGTCT CTATGCAGTA  
GAAGAAATGG CAGCGGAAGC AGCTCAATAC AAAGATCAAT TAGATCCAGC  
CACTGTAGAC TCGCTTACTC AAGTCTTTGG AAGCTATATT CTTGAGGTTG  
CACATAGAAA GCATGGCGGC TCTTACGTTT GGCTTGAATC TGAAAACTCA  
CCTGCCTTGG TAGTTGGTGA ACCAGAGTAC AGGCTAGCAC TCTCAACCTT  
CGCCAAGGTA CATGGCCGAC TTTCTGGCGA CGAAGCAGAT AATCTTATTT  
TCTTCTATCA AGGCTTTTCT GAAAGGCTTA AATCACCATC TCCCGGCATG  
AGCGCACTCT ACAAATGAAA CCCGAGTTTG GGGGCCCAAC AATGCGCTCA  
ACTGCCGCTC ACTTCGTTTCG CTGGACGTCC AAAAGCTACG CTTTTGGCCG  
CCCGTTAGCT TAATCGTTAT GCACAATAAA ACATGAAGAC AGCACTCATA  
TTTGTAGCTC TAATCTTTCT CTCTGGATGT GACAACTATC AGTCATGCCC  
TATAACTGGA AAATGGAAT CCAACGAAAA GCTAACTTTA GAAAGCATGA  
ATGAAACCGG CAGGATAACG GCAAAGCAAA GAGAGATTTT TGAGAACGGC  
TTCTTTGGAA AACTAGAATT AGACATAAAT TGCAGTAGCT TCACAACAAT  
ACTTGACGGC GTTACCGAAA CCTTTAATTA CGAGATAGTT CGCCAAACAA  
AAGATTCCGT CACCGTTAGC TATTACAGCA AAGCGCTGCA AAAACAAGTT  
GAGGTCACAT CTATTATCAA CGGAAATTGT TACTCGACAC CTATAGAGCA  
GTTAAATTTT AATGAGTATT TCTGCAGAGT CGAGTAGCGC ATAACAATTG  
ATTCAAGTCG TTCGCTTCGC TCACTGCGGG ACCGGCTAAA GCCGGCCCCT  
TAACCAAACG TTAGGCAAAG GCTCAATGGA TCCCATATTC CATAACATCC  
ATAGAAACGA CAAAGAGATT GAGGCGCTC ATCAACAATG CTCGAGCACA  
ATCAATCACT TCATTGAGAT GGTCAAAAAA GGGGGCGAGC CCACCTATAT  
GGCAAAGCTA CGTTTTCTTG ACCCTGACAA GTCTGAAAAA GAAGGTAAGA  
ATCATATTTT TTATTTGTGG TTATCTGAAG TGCTGTACCA CCCTGCAACA  
AATTTACTTT CTGGGGTATT TTTTGAAATC CCTGAAGGCT TTGAAAAGTG  
GCACCAAATA GGCCAGCGCC TAGGCTTTGA TCCAGAAGAT GTCTTTGATT  
GGATGGTAAT CGACAAAGGT CATGCTAAGG GTGCATACAC ACTAAAGGTA  
TCGCGAGAGC GCTTAACCAC CGAGCAAGAA AGAAAAGATT TTGACCGCTA  
TATTGGTGTG GCGTCATATG AGTAGCCTAA AATTAAGCGC TCACGCCTCA  
GCCTAACTAC TGTTTCAAGT CACTCGCTTC GCTCGTTCGG GACCGCGTTC  
CGCGGCCCTT TAACCAAACG TTAGGCGCAA GGGCAATATT GGTCTTCAGC  
ACCGAGTCAG GAAACACAAT CACCGAATCA GCGCGGTGTT CCTGAATCGA  
ATGGTCGCTG ACAGTTGAGG CCGTTATTTG TGGCCAGCAA AGGAGTTGCT  
TTCAGAGAAT GTGCACGTCA CAAATAACTT CCGGGGCCAA AACCGAAACG  
CCGTGCGCTC CGCCGGTTAA GCTCGGCGCT GGCATCATTT TCGGCGCTCG  
GCTGGGCAAT CTAACAATTG GCTCAAGTCG TTCGCTTCGC TCACTCGGGA  
CCGGCGAAGC CGGCCCTTGA GCCAAACGTT ATGCGAGCCA CCATGAATAG  
CGAAGAATTA TACAAAAAGC CTATGGAGTT AGAGTCCAAA TGCGAGCATA  
AAGCGGCAAT TTCAACTTAC AAAGAAATTG TTAAGAAATC TAACGATCCT  
CGCACTTCA TCGCAATTCG AGTTTGCCTC CAAAAATGTG GTCCTGGAA  
GCAATCCATC GAGGTATTAG AATCAGGAAT TGCCTGAAG CCTCACTATT  
GCGAGGGTGA TGCTCGTCTA TTTTTCAGCA AAGCACTTTT TAAATCAGGC  
AAAAAAGGCC TTGCGATAAA GCAATGGCAA CATGTATCAA AAATGCAACC  
TGAGTACCCA AGTTATGAGT CTGTGCAAAA TGAAGCCAAG AAAATGCTTG

CACAAAACGC	ATAACAATTG	GCTCAAGCCG	CTCGCTCCGC	TCACTCGGAC
GTCCGTAAGC	TACGCTTCCG	GCCGCCCTT	AGCCAAACGT	TAGGCACCAC
ATGCCCTCCA	TCAAGTCAGC	AAGCCAATAC	CAGCGCGCAA	CATCGCTCAT
CTTCTTAGTG	TCAGGCGCTG	CTTGGCTATT	CATCGTGCAG	TCTTCGTTGC
TGCCATTGAC	GGATGTCGCC	CGCCAAGAAA	TGGTTTGCCT	TAATATCGTT
CTTGGTATTG	CCTGTTTTGT	TATAGGTAGT	GCGGCAAAGC	GTCAGCGAGA
ATTCGCTG	CCTGACTGTG	GGAACGAAGT	AGATCAGAGC	TTACCTACAG
AGGGTATG	CGCCCCACTC	CTAAGGCTGT	GCAAGCACTG	CGATAATTCTA
TGGAATGTTG	GCAAGACCCC	AGACAGTTAA	AGTTACCGCC	TAACAAATTGG
TTCAAGTCGT	TCGCTTCGCT	CACCTCGGGAC	CGGCTAAAGC	CGGCCCTT
ACCAAACGTT	AGGCAACAGG	GGGTGACATG	ACGCAATGTC	CAAGGTGCCA
GCGCAATCTC	GCAGCTGACG	AGTTCTATGC	TGGCTCTAGC	AAAATGTGCA
AGGGTTGCAT	GACTTGGCAA	AACCTAAGCT	ACAACGCGAA	TAAGGAAGGT
CATGCCAACA	CCTTCACCAA	AGCGACATTT	TTGGCGTGGT	ACGGCTTATC
AGCACAGCGG	CATTGTGGGT	ATTGCGGTAT	ATCGGAGGCA	GGTTTTACAT
CCTTGCACAG	GACTAATCCA	CGCGGCTACC	ACATACAGTG	TTTGGGTGTT
GATCGCTCAG	ATTCGTTCCG	AGGCTATTCA	CCTCAAAACG	CTCGGCTCGC
CTGTTTTATA	TGCAACAGGA	TAAATCAAA	CATCTTCAGC	GCCAGTGAGA
TGGAAGTTCT	AGGTGAGGCC	ATTTCAAAG	CGTGGCATGG	TCGAGGAATT
GCCTAACTAT	TCAGTCAAGC	GGACGCAAA	CCCGCTGCGC	GGCTTTTGGC
CCGCTTATCT	CAAGCGTTAG	ATGAATAAAA	GCCTCCACAC	ATAGCCAGCT
TTACGGGAAC	GAAGTTGATG	CGACGCCTCT	TTACTGCCTT	AGTATTAATT
TTGATGATTT	CTGGGTGCTC	CTCCACATCA	AAAAC TGAA	GCCATAAACA
GCCGCCGAAC	AATTCAAGCG	ACACGACCGC	CATACTGAAA	TATATATTTA
CTGTTACCCA	TGGCATAGAA	GAGGGCTTTG	CCAAGTCACT	TGAGCAAGGA
AGTTACACCA	CAGATGAATA	TATGCTCATG	CAGCAGGCTT	TTAGCAATCT
TGATCTAAAC	AGATTAAACA	CTCTTTTATC	ACCCACTTTA	GATAAAAGCA
TGAACACCGC	AGACGTAAAA	CATTTTCATG	TTTTTATAAA	ATCTACTGCA
GGTAAGAACT	TGCTAACAGC	AGGAGAGTCT	AGCACTTCTT	TTTCCGGCAC
AATGGATAGA	GTTCGATCAC	TACCCTCCGA	ACAGCAGTCG	AAAATAAACG
AGTTCTTCCA	TGCCAGCTAC	ACAAAAAACA	CTTTAACAGC	CATGGGGGCT
CCAGAGGCGG	TACGCATTGT	TTATGCATTT	GGAGTGGAAT	CCATGTGCAA
TTACGCTATG	CGCAATAAAT	TCGAACTGTA	TATTTCTATT	ATAGAAAAAG
GCAAATGCCA	ATAAACCCACA	AAAAACAAAG	GCCGGATCGA	TCCATGTGAA
CATCAGCGTT	ACATCTAACA	TGTGGTTCAA	GCCGCTCGCT	TCGCTCACTC
GGGACCGGTT	AAAGCCGGCC	CCTTAACCAA	ACGTTAGAGG	ATTACATGCC
ATCACTCCAA	GAACTCCAAT	CGCCGATTGA	CTCAGCAATC	GTGAAGCTCA
TGATTGAGAG	CACTCCCGAA	ACATGAAGCC	AGATTATTCT	AACGTTGGTA
CGCGAATCTA	ATTCTCTCTG	TGTAGGTAAC	TTTACACATG	AGTTATCCAG
TCCTGAGGGG	CATGCACCAG	TTGGTCCGGC	AGAGAGCTTA	TTTGAAGACA
CTTACCAACT	CGATGAGTTA	TTCTACAGCC	ATGGTGAGCG	CTTATTCCAC
AAGGCAATTT	ATCGGGCTAA	TGCGGTTGGG	GACGGTTGGT	CATATCACGC
TGAGTTTGAA	TATGCGTAGC	ATCCCTCTAA	CAATTGGTTC	AAGCCGCTCG
CTTCGCTCG	TCGGGATCGG	CGAAGCCGGC	ACCTTAACCA	AACGTTAGAG
ATGGTCACTG	ATAAACGCGC	ACTAACCTTC	GGGCTACTCA	TAGCAATTCT
AGCTAGTATC	ATTTACACAAG	CGTGCTTTCA	TGGCCAAAAG	GTAACTCGCTT
CCGACGTGGC	TTTATATACT	CCATATTTCC	TATCTCCAAT	ATTTAGCATG
CTGATACAAG	CCACATCGAT	GCTGGCATGG	GCCATACCTG	GACTTTTATGT
AGGCTACTTA	TGCAAGAACA	AGCCAGCACA	ACATGGAGCA	AAATTGGGGG
CAGCATATGG	AATACTTCTT	GGATTAATTG	TATTCGCAAT	GCGAGCTTCG
ACCCAATTAA	CGTAGATTCT	AAGTTAATAA	TCGCAACATC	TGCTTTAACA
CAAAAAGCAA	AATATTCACT	GCACTTTGCG	CTAGTTGCTC	CTGCCGGCTA
TCTTCTTGCA	AAGCATCTGT	CAAACTCTCTA	ACAATTGGTT	CAGATCGTTC
GCTTCGCTCA	CTCGGGGACC	GGCTGAAGCC	GGCCCCTTAA	CCAAACGTTA
GGCAACTGAA	TGATCACCTG	CATTCCGGCA	CGTGAATTCC	TGCGTAAAGT
ATGCGGCCTG	TACGAAGCCT	CAACTAATGT	AGTTAAGTTG	CGTGTATGGG

CTTGTGGATA TGGCATCGCA ATGGATCTAA CTGTCAAAGG TAAATCTGTC  
CTTTGTGCGG TTGCGGGAGT ACTCCGCCAA GAGGTCTGAAT GCTTTGCTCA  
AATTGGGCTT CCGAACGTAA TTCAGTTAGT AGGCGACAAG GCGTCAGAGA  
ATCAACTAAA GCTCATAGGC ATGGAACCAC CAATCGAACT TCATATCTCT  
CGCGAACABA GCAGGCTCCA AGTTGTAATC TTGTACGAGG GTCAGGTAAA  
GGCTACATAT GTGCTGTCAG CCGCCTAACT ACTGGTTCAA GTCGTTGCT  
TCGCTCACTC GGGACCGGCT AAAGCCGGCC CCTTAACCAA ACGTTAGGCT  
TTCAATGAAA ACAGTTCAG TGAAAATATC AGAAGTCGAA CTAATAGAGA  
GTTTTGGGAA ATTCTGATC AATCAAGACT TAATCGACTA TGAAAATTCC  
CACTTCAGTG GCGACGACAA CCATAATGCA GATGTAGCCT TATCTTTAAA  
GCCAGGGAAA TGGCCAGGCA TTCAAGTCGA TAAACTACAC ATAGAAGTAA  
AGTCACACCA CTCAGAAGAC TCTCAAAACA CCATCAACAA AATATTCGGC  
CAATTACTAA AAGAAACCGG AAAGCGAAGC CTCGATAAAG AGAAAGAGTG  
CTTAGCTATA TTGTTCCCTT ACGAGCGCGG CGCATGGCCA GGTGGAACA  
ACAAAACAGT AACAAGAATT GAAGGTGAAG CTTATTACCG GAGGGGCTTT  
TCGAGAATCG ACAAACAGAC GTTTGTTAAA TTTGGTGAAT TGGTCGGTGC  
CAAATACATC CTTTCCTTTT CTACAGCATC AAACACATTG AACGTATTTG  
AATGGAAAAA TTTCTTAGAT GAGGAATTCA GCGCGATGAT CAGCCTAACA  
AATGGTTCAA GCCGTTGCT TCGCTCACTC GGGACCGGCT AAAGCCGGCC  
CCTTAACCAA ACGTTAGACG CACCGGAAAT TTTGCATGGG AAACCAGAAA  
ATGGATTGTC AGATAAACGA TACAAAGGTT GAGTGGGTTT CTCCAATACT  
GAAGCAATGG ATCAGCATCA ACAAAGAATA CGTCAAGCAA TATGATTTCA  
AAGACTGCCT CCACTGGTAT AACGAAAGGG CAAATATAAG TGCTTTTGCT  
GGTGCCGTTT GGAAGTCTGG AGTTTTGCG CTGGAAGAAT ATTCAACTAA  
AAAAGGCACC GAAGAAAACA GAGCCAATGG TCGTGTGAC CTATATTTCT  
CCAATGACAA CGAGCAAGCC ATTGTTGAAG CAAAAATGGA ATGGCTCTAC  
TTCGAAAGC GCACAAGACT AGATTTCAA GAAAAAATAG ATCGTGTAGT  
TGAAAAAGCA AAGAATGACA TAATTAACAG CCTGCATGCC AACCCCTACG  
ATCTAGGGCT TGGGCTTTCC TTTATTTGCA CATACTGGAA AAAGGGTTAT  
GACGCATCCG CCGACATGCA AGCCCTTAGA GCGCTTATGC AAAATTATAA  
CTGCGCATTT TATGCAATTT TTGAAAACAG CCCCAGACAAC GAAATTGTTA  
GCTCAAAGG CAATATCTGC AACGCTGTGA TTTTAGTTGG GACGGCGCAC  
AGCTGAATCG TGTGTGTGCG TCTAACAATG CGCTCAAAGC GCTCACTTCG  
TTCGCTGGGA TCGGCTAAAG CCGGCCCTT AGCTTAATCG TTAGCACTAG  
GACTTCCGAC CATCATGAGT GATAGAGACG AATTTTCTGC CCCAACAAAA  
AGAGCGCTAG CCGAAAGGAG TGGCTTTAGG TGTTCTTATC TTGGTTGCTC  
TAATGCAACC ATAGGGCCTA GTGAAGAATC AGAAACAGCC GTAGCAAGAA  
CGGGGGTGGC GTGTATATA ACTGCCGAG CCGCCGGCGG AAAAAGGTAT  
GACCCAACAT TAAGCCCTAC GGAACGAAGC TCAATCTCGA ATGGTATATG  
GATGTGCCAA ACGCATTGAG TTGAAATAGA TAGAGATGAG GCCCGATACA  
CATCGACCTT ATTAATCAC TGGAAAAATA TATCCGAGAG CCGAGCAGAT  
TATGCAAAAA ATCATGGCTG GGATATTTT GACAAATACC CCTTCCTTCA  
TATTGACTCG CTAGCCAACA TAGACCTGGC TCTTACCAA AGCCCTTCCT  
CAAATAGCCT TATCGGGAAT GCCATTACAG ACAGCTGCCT CCCTCAACTA  
TGGGGTAAAG AGCAATCTGT AATCATCAGA GACCTAATAA TAGAACTTTA  
TCGAAATGCC TTCGATCAG GCGAGGCTAG CTCAATCGAA ATATCCATAT  
CGGAGCAAAA ACTAGAAATA GTTACGATG GCAAAAAATT TGACATCTTC  
CAACTTCTTG ACCACCAGAA TGCAAACGGT GGCGCCGATA CCTTGCAAGA  
AATTGTAGAA AAATATGGCA GTAACCTTGT AGTCAACTAT AGCCACGAAG  
GCAACAATAA AATAATAATT CACAGGCTCT CTGACTTTTA CGCGCTTGCA  
CCATCCCTCC CGTGCGTAAT ATCACTGAGT GAATACGATG ACAAGGCCCT  
AGAGTTAGAC CTGGCTATTT ATGAGCGCTG CGGTGCCTG TACATAATTC  
TACCGTTGCA TTTTGTAGA TCAGATGTCA GGGGGCTAGA GTCGCAGCTA  
GCCGCTTTG AACCTAATGG AAAGCCAGTT TACATTGTAG GCTCAGATGT  
GGCAGAGCCT ACAAGAAAAG CAATTATAGA CAGGCTTCCC AACTTCACGT  
TCGTCCAAAA GCAATGCTAA CAATGCGCTT AACTGTCGCT CACTTCGTTT

GCTGGATAGT CAAAAGCTGC GCTTTTGTCT GTCCGTTAGC TTAATCGTTA  
GGCGCAAGGA GGGACCGTGA CTGAAACTGA GAAAATGGTG GGTAAGTTTCG  
TCAGCGGTTT TGGCGGGCAG AGATACCGAG AAATTTTGA AGTCCTCGAA  
TCCAGTAACC TTCGCCCCT GGGCAAGTCA AATACTGAAA CATTGCTATT  
TCAGCTTCGA GGGGCTGATA GTGAAATGCT AGATATTTTT GCCTTTCGCT  
TGGGGCCGCC GCCAGTAATT TCGTTTCCCA AATCATATTG GCTAGGTCGC  
CCCAGTGAAT TAAGCGCTCA TCTATCCAAT TTTTCATTCT CGGAAAAGCC  
AGCCATAACA GGCCCGGTTT CTGACTCACA GTATTCGGCA GGCCAGGTGG  
AAATCACCCG CTCTACTCAT GAGAGGATTA TTGAGGTTTG CAACCGTGTC  
TGTGCTTCCC TGCAATAAGC GCCTAACAAC TGGTTCAAAT CGCTCGCTCC  
GCTCGCTGGG ACCGGCATAG CCGGCCCTT AACCAGCGT TAGATGCAAA  
TAACCTGAGG GGCACATGCA AGACTTTGGG TCCAGACGAA ATGCATCATT  
AGAGGACAGG GCTGCGGCTG AGTCTGTTAT TGAACGTGTT TATCTTGCGA  
TACAGCAGCT TTGCACAGAG ACTGGTGACG TAAGAAATCG GCTTCAAATA  
GCCGTTATGA CTCTATTGCC CCTTCAGGCG CGTAACTTCC CCATTGCGTT  
GCAGCAAGAC TTCGATTGGA TTGTGAGAGA ATCAACCAAA TACAAATCAC  
CATATCCGCA GTTTCGGGGC GACCTTGAAG CAACGATGAT GCGAATAAGG  
AACTCAACTG GGCAAAAAAT CGCGCAAAGA ATTTTCAATA TTTACTCGTC  
GCTACAAGAC ATTCGAGGTT TTCCCCTGCT TGAATACAGG GCAATAGATG  
AGTAAGCATC TAACAATTGG TTAACCCGT TCGCTTCGCT CACTGGGACC  
GGCTAAAGCC GGCCCTTAA CCAAACGTTA GGTAACCAAG GGAAATTAC  
TTGAGTTGTT ATGTATTGGG CACAAACAAC CGCCATTAAA GGACGGTTTT  
ATAGTAAATT TCATCGGACT GTTGAACATA AATGCTTATA CGCTTTGCTC  
TACTACTTGC TGTTATGCTC TCGCTGCAT GCTCGTCAA GCAAAATCCA  
ACGCCGAAGT GTACTGCCAG CGTCCCCCG CCCTCTTAC CCGAAACATC  
CACAGTATGC CTAGGGGAAA GATGTAATTG GGAGGTGCTA TTTCCGTCAG  
GAAAATACCC TGCATCCACA GAAGGCTGCA GGGCGCCTGT GGTGCAGAAC  
CAGCCTTCTT CCTACCCGCG AGAAGCACTT GATCAGTGA TTGAGGGGTA  
CGCTTGGGTA GCGGTTTTTC TGAATGCCGA CGGGGTCCAA ACATCAGCAA  
AGGTACTTCA ATCATCGAAT AAAATTTTCG ACAGAAATGC CTTGCTACAG  
GCCAGTAATA TATTTTTTGA GCCTATGAAA TGTCATCCG AGCGTTATGA  
TTCCGTTGTT CTGATGCCAT TAACTACCG CATACTCCCC TAGTAGCGGG  
ATTGATCCTT ACAAATTC AACTTACGT CTACTTACGT CCAAGTTGAA GTAGGCAGTT  
TAACAACCTG TTCAAGCCGC TCGCTTCGCT CACTCGGGAC CGGCTAAATT  
CGGCCCCCTT GGCAAACGTT AACTATCAGA AGGGCGGTTG ATGTCAAGAT  
TTGCGCTCGC GTTGATTAC GGAGTACCAA CGGGTTTTCT TGTCATTTGT  
ACTTTGTTTG TCTGTTTCAT CTACCTCAAC CGATTCGAGA AAGTTGGAGG  
ATACTCAGAC GGGTGGGGTT TTGTTGGAAG AGTTGTCTGC GCATCTATAG  
CTATGGTTTT CGTGTCCGCA GTTGGCCATC TTCTTATTGA AGCGGCAGTC  
AACTGGGGGC TGCAGCAGCT TGGTTATGAG CTGCCAACT ATGAAAAAG  
AAGGACTTGT AGTAGCTGCA AGCCGAGCAC TCCAGGTGAC TACATGTTTCG  
GCTTGCTCCT CGGGGGTGTG CTTGGCGCCG GCTCGGCAAT TTGGCTCTGG  
ACGCGCCTGG CGCTCCGATA TGCGCTGTTT CGCGGCGAAA ACTGATAGCT  
GAACCTTCCA TCGAGGAGAT GCAAAAGCGC TGCTGCGCGC CATCTACAAA  
GACCCGAAGC ACCTCATCCA GCGGCTCTCA GCCCGAGCCT GACTGGCTGT  
GGCTATCAAC ACCTCTTCGA TACCACTACC CGCCAGAAAC GACAAAGCCC  
TGCAAAAAGC AGGGCTTTGT CTTTGGGGAT CTGGAGCGGG CGAAGGGAAT  
CGAACCCTCG TCATGAGCTT GGAAGCTCA GGTAATGCCA TTATACGACG  
CCCGCTCGGG CGGCTGACTT TTTACCAGAA TCGCCCGGGA AGGTGAAGCC  
GGGCGCGCGT CTTGCGCCCC TTTTATTGCC GGGCGCTTCA TAGCGCCACG  
GCCCCTGGCT CTCGTTCCAC GCTGCGTGCG TGGCCCTGCG TGGGTGCCAG  
CAGGAAGGCC AGCAGGGCAT CGCGGCTCTG CATCCAGGCG GCCTTGTGTT  
CCATGTCGAG GAAGTGGCCG GCCTGGGCGA TGGTGCGGAA CTCGCAAGTG  
CGCACGTACT GGGTGAACAG GCGCGCGTCA GCGGGGTGG TGTACTCGTC  
CCACTCGCCG TTGACGAACA GCAGCGGTAT CTCGATCTGC CCGGCGAAGC  
TGACGCAGGA GCGCCCGCCG TTGTTGAGCA CGGTTTCCAC GTGGTGACTC

ATTTGCTCAT ATTCATAGCG CTCAGGCCG GTGACGTGTC GATGGTTGTA  
 GCGCTTGAAC AGCGAGGGCA GGTGCTTGCC GATGGTGCCG TTGAGCACCA  
 TGCCGATGCT CTCGCGGTCG CACTCGCGCA TCACCACCAG GCCGGCGCGC  
 AGGTAGCCGA GCATGGCGCT GTTGACGATC GGCAGAGAAG AGTTGATCAC  
 CGCAGCTCG ATCCGCGATG GACGCGGGC CAGCGCTTG AGGGTGCGCA  
 TGCCGCCCCA GGAGAAAGGA CAGCAGCTG TTCGACGGC AAAATGTTG  
 ACCAGTCCA GGAAGATGTC GGCTTCTTC CTCGCSGCTG AAG



Figure 3B-1  
SEQ ID NO:2

AAGCTTCTGG TACGAACCTG GGGGCGCTCC GGCACGCACA AGGGCATCGA  
CATCTTCGCC CGCCAGGGCA CCCCAGGTGCT CGCCCCCAGC TACGGCATCG  
TGGTGTTTCG CGACGAGCTC GACATGGGCG GCAAGGTACT GCTGATGCTC  
GGCCCCAAAT GGCCTGCA CTACTTCGCC CACCTCGACA GCTACAGCGC  
CCTGCCCCGC CAACCCGTAC TTCCCGGCGC CCCACTCGGC ACGGTAGGCA  
GCACCCGGCAA CGCCCAGGGC AAGCCGCCCC ATCTGCACTA CTCGATCGTC  
ACCCTGTTGC CCTATCCCTG GCGCTGGGAC AACAGCACTC AGGGCTGGAA  
GAAAATGTTT TACCTCGACC CCACGCCAAT GCTGAACGAA GCGGCAGTAG  
ACAGCCGAAA AACCAGCCAG TAGCGTCGCA GGGGAATGCA CCACCGGTCT  
TGCCCGATCC GCCTGTCCTT TTACCAATCG CAGAAGAGTC GCTTTTGTCTG  
AATCGCCTGT GAGGAAAAAC AAGGACTTGC TGGACGACAA GGAACGTTAT  
GCGACACAAG TGGCGGAATA AATTACGCCA TTTGTGTCTG CTACTTATAG  
TTATATGCTG ATCTAGATAT GAAGTACAAA AACATAAAAT CAGCAATCCA  
CAATTTCGGG CACAGCTTTG TAAGCTCAGT GAACTATGTT GACCATGATT  
TCGTTGCCGA CGAAATTGGG AAGATTCACA AGAAAGGCTA TGATATTGAA  
ATAAACTGGC TTACAAGGGA GTTCAAGCCC GCTCAGCTTG AGTCAGAGAG  
AATAAAAAAA TCAATTGGTT ATTGGGGTGA CAACCTAAAG AAACATTGTG  
CATCCCATAG CGTAAATCTG GAAAATCTAT GTTCTTTATC GTTTATCTGG  
CCGACAGGTC AAAGTAAATA CATGCATGCC ATTGACGACA AAGGCACAGA  
ACACAAAATT TACATCAATG AAGCGCAGTG ATACGCATAT AACAATTGGT  
TCAAACCGTT CGCTGCGCTC ACTGGGACGG GCTAAAGCCC GCCCCTTAAC  
CAAACGTTAT GCGAGCCACC ATGAATAGCG AAGAATTATA CAAAAAGGCT  
ATGGAGTTAG AGTCCAAATG CGAGCATAAA GCGGCAATTT CAACTTACAA  
AGAAATTGTT AAGAAATCTA ACGATCCTCG ACACCTCATC GCATTCCGAG  
TTTGCCCTCA AAAATGTGGT CACTGGAAGC AATCCATCGA GGTATTAGAA  
TCAGGAATTG CACTGAAGCC TCACTATTGC GAGGGTGATG CTCGTCTATT  
TTTAGCAAAA GCACTTTTTA AATCAGGCAA AAAAGGCCTT GCGATAAAGC  
AATGGCAACA TGTATCAAAA ATGCAACCTG AGTACCCAAG TTATGAGTCT  
GTGCAAAATG AAGCCAAGAA AATGCTTGCA CAAAACGCAT AACAATTGGC  
TCAAGCCGCT CGCTCCGCTC ACTCGGACGT CCGTAAGCTA CGCTTCCGGC  
CGCCCCCTAG CCAAACGTTA GGGGCCAAGA TGGATCTTCG CCAGACAAAG  
CCAATACTAG TTACAGTCTT AGCCACTGCC TTGGTGCCAT TGGTTTTTGG  
CTGGTATGCG TATTGGGAAA ATCCTCAAGG CATACTTTTG TACACTCCGG  
TGGCCGGCCA TCCCCATCCT CAGGGCTCTC CAGCATTTCC TATTGGAGTA  
ATGGTTGGGC TGGCCGCTTC ATTTCTGCTC TCTTTGCTTT TTGTAGGCCT  
AGGGGGAATC GCTGCATACA TAGCAAGTTC AGTGAGCTCA AAGGCTAGGG  
CTAAGCTGTT TTGCAAAATC GCAGTCACAT CCCTGGCTAC TTCAACTATA  
GGAGCTGCAG TCTATGCAAT GCTCCCCTAA CAAATGGTTC AAAGCCGTTT  
GCTTCGCTCA CTCGGGACCG GCTAAAGCCG GCCCCTTAAC CAAACGTTAG  
GCAGCACATA TGAATCGTTC GTGCCTATAC ATGTTTATCG CCTCAGCCTT  
GATAGCGTGC GGCATCCAC CTCTATTGGT TACGCCACTG CCAAATGGCT  
ACAATTTCCA TTCCAACGGC GGGGAGTTTG GCTACATCAA GAATCCAGAT  
GGATTAAGGC TCGCCGAGTA CTTTGGTATT CGTAATGATG GTCGCGAAAC  
CTGGTGCACT GACTTTTCAT GGGAAAGCGA TATCGTCATT TGTAAGCTTA  
TTGAATATAG CCAGCATGGA TTTGACGCAT CGCATAACAGA GTTTTCTGTA  
CTTGACACAA AACTAGCGA GGTAGGGTA TTTCCCGATC AAGCGTCTGC  
TCAAATTTT TGGGCCGCAC GCTTTAATTC AGGACTACCT CAGCTTCACC  
GGCACTACCC TTCAACCTCA GAGAAGTAAT ATTTTGTGTG TCAGTGCAGC  
CTAACAATGC GCTCAACTGC CGCTCACTTC GTTCGCTGGA CAGTCAAAAG  
CTGCGCTTTT GCCTGCCCCT TAGCTTAATC GTTAGAGGCT TATTTAGCTC  
ATGCGCATAG ACATAGACTT TTCAATATTC ACGCTCGCAC CGTCGACCGA  
AGGCGTAATA TCAGGAAAAA TCGAGGTCAG TGAACCTACCT AGAACTGGCG

[illegible]

ACTTTTAAACG GCTCAAAGAG CATGGATCGC GTTCAGGGAA GCAAACGTGTG  
CCACTCAGTA CGAAATGCAC AGATCTGGCA CTATTTCGCAA CAGCATCTAT  
CTAGCCTGCA AAGAAAAGCG TGCCAAGCAG CGAATAAAGG GAGCTTCAAA  
ATTATGCTCC GTACTAGCCC TAACAAATGG TTCAAGTCGT TCCGCTTCGC  
TCACTGCGGG ACCGGCTAAT GCCGGCCCCT TAACCAAACG TTAGGCCGAC  
AATCGCAATT CCTAGGACTG CACGTGAACT GGATCCGCAA AATGTTTCGG  
CGCACAGCAC TAGCGCCGCC CCAACATCGC GAGGACGAAG CTGTCAGTAC  
AAGCCAAGAA GGAACGCCCT CCTTTCGTCA TTTGACAGTT GAGAATTCAT  
GGGGAAGTTG AGGGCGGAGC TATTCCTTCA GTCACCACCC GAGAACATCC  
TCAGAAGATC TGTTTCTCGT TTGGCGTGCC TAAGTTCGGA TGGTCAACGT  
TCGAGATCCA TTTCGTCGGA AATGGCCACT TCATCTGCGG CATCTCTGAC  
ACTCCAAATG ACTTCTACGG TGACTTGGCT ATCGCCCTGG CTGAGCAGAA  
AAGTTCCTTT TCGGTAGCGG CGCACCTTGA GCCTGAGACC TTTGCCTTCT  
ACATCGTTGA TTCGACAATG TACTTGTGCA AGTTCGATGA ATTCGACGAT  
TATGAGTCCG CCGCCGAAAG CCACGAACAG TTGGTCTCCC ACAGCTTTAT  
GTCCATTGAA GTATCTAGGG AGTACTTTCA GAAGTCTCTC AGGACCTTGG  
CCGTCCAATG GCCGGATACG CCTTCAAGAG ACTGGGCGCA CCCATTTCCA  
CGTGCGCAGA TTGAAGGCTG ACTGCCTAAC TATTCGCTCA AGCGGGCAGC  
GTTAGGCGCC CTCATTCGGA GTCACGCTAT GGCAACCCGA GAAGAAACAG  
AAGTAGCCAT TGCTGCTCTT CGCAGCGAAC TCAATGGCAA CGAATCGGAA  
TACAGCTTTC ACATTCCCCG TTGGGCGCCA GAAACATCAG TCATGGGATT  
TCGCTGGATG CAAAGCCAAC TGTGGGAAGG CTTCTACGTA AGCTATCGCG  
TAGAGCACTC GGCCAAGCGC GTCGAATTCA AGTGTCTGGG GTACGGCGAG  
CCCGAGCCGT CTTGGCTGCA AGTTGGCTAG GGGGCCGGCA AGATGCAATC  
GCGGCGAGCG CCTAACACTG CAGTCAACCG GACACCAAAC TGTACGCAGT  
TTGGTTCCCT CCGCTGCGCT CCGGTGCCGG TTACTTTCAA CGTTAGGCAA  
CTCAGATGAG TGCTCCAGAC GCAGAACTTC TCGCATTGTT AGCCTACCGA  
ATGGAAGCTA TTTCCATTGG GCATTTGGCA TTACGCCATC ACATGACGTG  
GGACGAAACA CCTTCAATGG AGGTGTACTT CAATGGCATA CAAGTACTCG  
AGGGAAAGGC CACGGGTTTC ACTAATGCAG CCATTGAGTC CGCAATTATT  
CATTGCAGGG CAATCCTTGG AGTTTGTTGG GCTGCAGTCC TCCAGACACT  
CTTCCACAGA AATTGCAGAG CGCACTCGAC GCAACAATCC CGATGACTAT  
GGCATTGAAA GCTTCAATGG CTTATCAATG CTAACCAAGG AAAAAGCACT  
AGCCTACTAC TCTGGCGAGC TGCCAGAAGC GGAAGTTGCT CTAGCGCTCA  
TATTCCACTC AGCGAACAAA GGGCTTGAC ACACCTACAGT GTCTTTTACG  
CGTGACAGTG GCGACGCCCA CCTGATGGAA ATTGCATTTT GCATCGTACC  
AATCCTGCTT GTAAATGGCT TCTACGCTCC ACTGGAAATC ACGCCACCAA  
AATATGAACT GATTTACGCG CCAAGAGTCG CCATAACAAA TGTTCAAGT

Figure 3C-1  
SEQ ID NO:3

AAGCTTCTGG TACGAACCTG GGGGCGCTCC GGCACGCACA AGGGCATCGA  
CATCTTCGCC CGCCAGGGCA CCCCAGGTGCT CGCCCCCAGC TACGGCATCG  
TGGTGTTCG CGACGAGCTC GACATGGGCG GCAAGGTACT GCTGATGCTC  
GGCCCCAAAT GGCGCCTGCA CTAATTTCGCC CACCTCGACA GCTACAGCGC  
CCTGCCCGGC CAACCCGTAC TTCCCGGCGC CCCACTCGGC ACGGTAGGCA  
GCACCGGCAA CGCCAGGGC AAGCCGCCCC ATCTGCACTA CTCGATCGTC  
ACCCTGTTGC CCTATCCCTG GCGCTGGGAC AACAGCACTC AGGGCTGGAA  
GAAAAATGTT TACCTCGACC CCACGCCAAT GCTGAACGAA GCGGCAGTAG  
ACAGCCGAAA AACCAGCCAG TAGCGTCGCA GGGGAATGCA CCACCGGTCT  
TGCCCGATCC GCCTGTCCTT TTACCAATCG CAGAAGAGTC GCTTTTGTCTG  
AATCGCCTGT GAGGAAAAAC AAGGACTTGC TGGACGACAA GGAACGTTAT  
GCGACACAAG TTGCGGAATA AATTACGCCA TTTGTGTCGT CTAATTATAG  
TTATATGCTG ATCTAGATAT GAAGTACAAA AACATAAAAT CAGCAATCCA  
CAATTTTCGGG CACAGCTTTG TAAGCTCAGT GAACCTATGTT GACCATGATT  
TCGTTGCCGA CGAAATTGGG AAGATTCACA AGAAAGGCTA TGATATTGAA  
ATAAACTGGC TTACAAGGGA GTTCAAGCCC GCTCAGCTTG AGTCAGAGAG  
AATAAAAAAA TCAATTGGTT ATTGGGGTGA CAACCTAAAG AAACATTGTG  
CATCCCATAG CGTAAATCTG GAAAATCTAT GTTCTTTATC GTTTATCTGG  
CCGACAGGTC AAAGTAAATA CATGCATGCC ATTGACGACA AAGGCACAGA  
ACACAAAATT TACATCAATG AAGCGCAGTG ATACGCATAT AACCAATTGGT  
TCAAACCGTT CGCTGCGCTC ACTGGGACGG GCTAAAGCCC GCCCCTTAAC  
CAAACGTTAT GCGAGCCACC ATGAATAGCG AAGAATTATA CAAAAAGGCT  
ATGGAGTTAG AGTCCAAATG CGAGCATAAA GCGGCAATTT CAACCTACAA  
AGAAATTGTT AAGAAATCTA ACGATCCTCG ACACTTCATC GCATTTCGGAG  
TTTGCCTCCA AAAATGTGGT CACTGGAAGC AATCCATCGA GGTATTAGAA  
TCAGGAATTG CACTGAAGCC TCACTATTGC GAGGGTGATG CTCGTCTATT  
TTTAGCAAAA GCACTTTTTA AATCAGGCAA AAAAGGCCTT GCGATAAAGC  
AATGGCAACA TGTATCAAAA ATGCAACCTG AGTACCCAAG TTATGAGTCT  
GTGCAAAATG AAGCCAAGAA AATGCTTGCA CAAAACGCAT AACCAATTGGC  
TCAAGCCGCT CGCTCCGCTC ACTCGGACGT CCGTAAGCTA CGCTTCCGGC  
CGCCCTTAG CCAAACGTTA GGGGCCAAGA TGGATCTTCG CCAGACAAAG  
CCAATACTAG TTACAGTCTT AGCCACTGCC TTGGTGCCAT TGGTTTTTGG  
CTGGTATGCG TATTGGGAAA ATCCTCAAGG CATACTTTTG TACACTCCGG  
TGGCCGGCCA TCCCCATCCT CAGGGCTCTC CAGCATTTCC TATTGGAGTA  
ATGGTTGGGC TGGCCGCTTC ATTTCTGCTC TCTTTGCTTT TTGTAGGCCT  
AGGGGGAATC GCTGCATACA TAGCAAGTTC AGTGAGCTCA AAGGCTAGGG  
CTAAGCTGTT TTGCAAAATC GCAGTCACAT CCCTGGCTAC TTCAACTATA  
GGAGCTGCAG TCTATGCAAT GCTCCCCTAA CAAATGGTTC AAAGCCGTTT  
GCTTCGCTCA CTCGGGACCG GCTAAAGCCG GCCCCTTAAC CAAACGTTAG  
GCAGCACATA TGACTCGTTC GTGCCTATAC ATGTTTATCG CCTCAGCCTT  
GATAGCGTGC GGCGATCCAC CTCTATTGGT TACGCCACTG CCAAATGGCT  
ACAATTTCCA TTCCAACGGC GGGGAGTTTG GCTACATCAA GAATCCAGAT  
GGATTAAGGC TCGCCGAGTA CTTTGGTATT CGTAATGATG GTCGCGAAAC  
CTGGTGCAT GACTTTTTCAT GGGAAAGCGA TATCGTCATT TGTAAGCTTA  
TTGAATATAG CCAGCATGGA TTTGACGCAT CGCATACAGA GTTTTCTGTA  
CTTGACACAA AAAGTAGCGA GGTTAGGGTA TTTCCCGATC AAGCGTCTGC  
TCAAAATTTT TGGGCCGCAC GCTTTAATTC AGGACTACCT CAGCTTCACC  
GGCACTACCC TTCAACCTCA GAGAAGTAAT ATTTTGTGTG TCAGTGCAGC  
CTAACAATGC GCTCAACTGC CGCTCACTTC GTTCGCTGGA CAGTCAAAAG  
CTGCGCTTTT GCCTGCCCGT TAGCTTAATC GTTAGAGGCT TATTTAGCTC  
ATGCGCATAG ACATAGACTT TTCAATATTC ACGCTCGCAC CGTCGACCGA  
AGGCGTAATA TCAGGAAAAA TCGAGGTCAG TGAACCTACCT AGAACTGGCG

AGATAATTTT ATTCTCCTTT GCGCCAAACA AGTCTAAATT CCCGGCAGAG  
 CCAAGATTCA ACCCGTTGCT TAAAGTTGAG AGAGTGATTG ATAGCGTAAA  
 TGGTCAGAGT CCAGCTCTTC AGTTAGAGAA TCTGATGCTA CCAAACAGAG  
 AAAGTGTGCG TGAAGTCACT GCTTTCCTAG AGCAAGGCTT TGGCCTATTT  
 TTCAGGCCAA CCGGTGAGTA ATCCTCTAAC AATGCGCTCA AATCGCTCAC  
 TACGTTGCGT GGGACCGGCT AAAGCCGGCC CCTTAGCTTA ATCGTTAGAG  
 GTCAGCACAT GGCAGTGCAG CAACTCGGGC CAACCACAGT ATCCGTAACC  
 GAATTTGTCAT GGGACGGAAG CGATCTTGGA AATACTGAGG CCAATGAATT  
 CTGGTCACAG CTCTCTGCTC AGCTTCAAAA AATAGCTATC TCTGAGTTTT  
 TAGCTGGCAA TCGCCCCAGC AGCATTCTTC GCAACGACCC ACGAAACATT  
 ATTGTTCTCT CATTTTCGGC GCGCCAAAG TTCATTAAAA TCAACCACTG  
 GCTCTCTGCG TGTACACACA GAATTTCAAC ACGGAAATTA CTGCTACGAC  
 GGAAACGGCC TGTACTTACG AAAATTTAGA GTCTGGCGAC TTTCTTGCAT  
 TCGACACAGC GCGGTTGGTG CATGCCCTCT AACAAATTGGT TCAAGTCGTT  
 CGCTTCGCTC ACTGCGGGAC CGGCTAAAGC CGGCCCTTA ACCAAACGTT  
 AGGGCACCCG GCATGAGAAA TGAAGACGGA ACCTTTTGCA AAGACTGCCA  
 CCATCAACTT GATGAAACAC TAGCATCTAG CGCAAATTAC TCATGCCCCA  
 ACTGCGGCTC CACAAAAAAA TACATGAACA TGTCCATCAC TGATGGAATT  
 GGCCTATACG ACTCTTTGGG TGCCCAAGCT AAAGATCCAA GTTACCCGGC  
 AAAAGAAAAA TCAGATGGGA AACATTTGTT GGCTGGGAAC GCAGTCATAA  
 ACTGCAAAAA ATGGTTTACA AGACAAGAAC TATCGATCGA ACCAATGACG  
 CATACCAAGA AATAGTAGTC GACCTTAAAA CAGGGGGAAT AATTCATCAC  
 TGTGAAGAGC CACTTTCAGA GCAYTTKGGC CATGGCACCG CAAAACCAAA  
 GCCCTAACAA ATGGTTCAAG TCACTCGCTT CGCTCGTTCG GGACCGGCTA  
 AAGCCGGCCC CTTAACCAAA CGTTAGAGGT TACCTGTGAC AGATTGCGC  
 CCGTTACTGA TCCCTGCCTC GCAATATGAT ACGAGCGTTC TTCTCGCCGA  
 ATGGCAATGG CTCACCCCCA AAACGGATAC GCCACTTTTT ATTTCCATAT  
 TCGGAGACTG GGTATTTGGC AACCCTAATG GAAGTTTGTG GGTCTTTCA  
 CTCCTAAAAG GCACTTACGA GCAAGTAGCC GCAAACCTA ACGAGTACAA  
 CACCTCAAC AAATCGGCGG AGTGGATTGA TCAAACATTC ATCGCCAGTT  
 GGCAGTCTAT TGCCGCAGGC CATGGGTTAA TCCCAGAACC AAACCAATGC  
 CTCGGCTGGA AGGTTACCCC ATTATTAGGT GGAAGTTTGT AGCCAGCCAA  
 TCTCCAACCTC TTCAACATGT CGGTGTATCA ATCGCTTATG GGTCAACTTC  
 ATCGACAGCT TAGCCAAAAA CAAACCCCGG CAAGTAAAAA ACCATGGTTC  
 CAGTTCTGGT AACCTCTAAC AACTGGTTCA AGTCGTTTCG TTCGCTCACT  
 GCGGGACCGG CTAAAGCCGG CCCCTTAACC AAACGTTAGG CGCAAGGGCA  
 ATATTGGTTA TTCAGCACCG AGCCAGGGAA CACAATCACC GCATCAGCGC  
 AGTGTTCTTG AATCGAATGG TCGCCTGACA GTAGAGGCCG TTATTTGTGG  
 CCAGCAAAGG AGTTGCTTTC AAAGAATGTA CACGTACAA ATAACCTCCG  
 GGGCCAAAAC CGACACGCCG TGCGCACCGT CGGTCAAGCG CAGCGCTGGC  
 CTCATTGCA GCGTACGGCT GGGCAATCTA ACAATTGGCT CAAGTCGTTT  
 GCTTCGCTCA CTCGGGACGT CCAATAGCTG CGCTATTGGC CGCCCTTAG  
 CCAAACGTTA GGCCAACATA CTCAACGCAT GAAAACAAAA TATCACATAA  
 ATATAATTAT ATTTCTCGAA ATCATAATTC CTTTAGCACC AATAATTTGG  
 GCAATTTTCA CTCAGTCAAG CCCC GGCTTT GGCCCAACCC TTATATCAAT  
 GTCATCCTG CACATCGTCG GACGAATAAT TAGCCGAAGC ATCCCTGCCA  
 GCTGTGACTC ATGTGCTGAA AAAATAAAAC CCAAAGGAAC CTCGCAATC  
 TACTACAACT GTCAAAAGTG TGGATTTAAA TACTCAAAAA CACTTAACAG  
 CAGCAAAAAC TTCCATAACC ACTAACCAGA AAATCACTAA GCGGCCATCA  
 GTTTATAAGC GCCGTAAGCA CTAAAGACTT GTACAAGCCT AACAACTGGT  
 TCAAGTCGTT CGCTTCGCTC ACTGCGGGAC CGGCTAAAGC CGGCCCTTA  
 ACCAAACGTT AGGGCACTCA ATGCATCGCT TCCTAGCCAC ATGCCTACTA  
 GCTACATCTA TTAAGGCATA CGCAGAACCT GAAAATAATA TCGACTGCAG  
 CAACGCATTC TCAACGCCGG ACATTGAACA TTGCGCATCA ATCTCTCTTG  
 AGAAAACAGA GAAAGAGCTA AATTTAGCAT ATCAAAAATT AGTCAAAGAC  
 CTTTCTCAGC CAAACAATGA ATACGAAAAT TTCACCGAGT ACAGGAAAAA

ACTTTTAACG GCTCAAAGAG CATGGATCGC GTTCAGGGAA GCAAACGTGTG  
CCACTCAGTA CGAAATGCAC AGATCTGGCA CTATTCGCAA CAGCATCTAT  
CTAGCCTGCA AAGAAAAGCG TGCCAAGCAG CGAATAAAGG GAGCTTCAAA  
ATTATGCTCC GTACTAGCCC TAACAAATGG TTCAAGTCGT TCCGCTTCGC  
TCACTGCGGG ACCGGCTAAT GCCGGCCCCCT TAACCAAACG TTAGGCCGAC  
AATCGCAATT CCTAGGACTG CACGTGAACT GGATCCGCAA AATGTTTCGG  
CGCACAGCAC TAGCGCCGCC CCAACATCGC GAGGACGAAG CTGTCAGTAC  
AAGCCAAGAA GGAACGCCTC CCTTTCGTCA TTTGACAGTT GAGAATTCAT  
GGGGAAGTTG AGGGCGGAGC TATTCCTTCA GTCACCACCC GAGAACATCC  
TCAGAAGATC TGTTTCTCGT TTGGCGTGCC TAAGTTCGGA TGGTCAACGT  
TCGAGATCCA TTTCGTGCGA AATGGCCACT TCATCTGCGG CATCTCTGAC  
ACTCCAAATG ACTTCTACGG TGACTTGGCT ATCGCCCTGG CTGAGCAGAA  
AAGTTCCTTT TCGGTAGCGG CGCACCTTGA GCCTGAGACC TTTGCCTTCT  
ACATCGTTGA TTCGACAATG TACTTGTGCA AGTTCGATGA ATTCGACGAT  
TATGAGTCCG CCGCCGAAAG CCACGAACAG TTGGTCTCCC ACAGCTTTAT  
GTCCATTGAA GTATCTAGGG AGTACTTTCA GAAGTCTCTC AGGACCTTGG  
CCGTCCAATG GCCGGATACG CCTTCAAGAG ACTGGGCGCA CCCATTTCCA  
CGTGCGCAGA TTGAAGGCTG ACTGCCTAAC TATTCGCTCA AGCGGGCAGC  
GTTAGGCGCC CTCATTGCGA GTCACGCTAT GGCAACCCGA GAAGAAACAG  
AAGTAGCCAT TGCTGCTCTT CGCAGCGAAC TCAATGGCAA CGAATCGGAA  
TACAGCTTTC ACATTCCCAG TTGGGCGCCA GAAACATCAG TCATGGGATT  
TCGCTGGATG CAAAGCCAAC TGTGGGAAGG CTTCTACGTA AGCTATCGCG  
TAGAGCACTC GGCCAAGCGC GTCGAATTCA AGTGCTGGGA GTACGGCGAG  
CCCGAGCCGT CTTGGCTGCA AGTTGGCTAG GGGGCCGGCA AGATGCAATC  
GCGGCGAGCG CCTAACACTG CAGTCAACCG GACACCAAAC TGTACGCAGT  
TTGGTTCCCT CCGCTGCGCT CCGGTGCCGG TTACTTTCAA CGTTAGGCAA  
CTCAGATGAG TGCTCCAGAC GCAGAACTTC TCGCATTGTT AGCCTACCGA  
ATGGAAGCTA TTTCCATTGG GCATTTGGCA TTACGCCATC ACATGACGTG  
GGACGAAACA CCTTCAATGG AGGTGTACTT CAATGGCATA CAAGTACTCG  
AGGGAAAGGC CACGGGTTTC ACTAATGCAG CCATTGAGTC CGCAATTATT  
CATTGCAGGG CAATCCTTGG AGTTTGTTGG GCTGCAGTCC TCCAGACACT  
CTTCCACAGA AATTGCAGAG CGCACTCGAC GCAACAATCC CGATGACTAT  
GGCATTGAAA GCTTCAATGG CTTATCAATG CTAACCAAGG AAAAAGCACT  
AGCCTACTAC TCTGGCGAGC TGCCAGAAGC GGAAGTTGCT CTAGCGCTCA  
TATTCCACTC AGCGAACAAA GGGCTTGAC ACACTACAGT GTCCTTTACG  
CGTGACAGTG GCGACGCCCA CCTGATGGAA ATTGCATTTT GCATCGTACC  
AATCCTGCTT GTAAATGGCT TCTACGCTCC ACTGGAAATC ACGCCACCAA  
AATATGAACT GATTTACGC CCAAGAGTCG CCATAACAAA TGTTTCAAGT

Figure 3D-1  
SEQ ID NO:4

GCCGGCCTGC AACAGAGCTT CAAGGCCGCC GGTGTCGGCA TTCGGCGACC  
GCGTGCCGCC TTTATGCCAA CACCTGGCAG ACGGTGGTCG GCATGTTTCG  
CACCGCCAGC TTGGCGCCAA CTGGTCGTCC TGCTCGCCGG ACTTCGGCAC  
CCAGGGCGTG ATCGACCGTT TCGGCCAGAT CGAACCCAAG GTGCTGATCG  
CCGCCGCCGG CTACCGCTAC GCCGGCAAGA ACCTCGATCT GACCGCCAAG  
CTCAACGAAA TCCTCGAACG CCTGCCCTCG CTGCAGCAAC TGGTGGTGGT  
GCCCTACTCC AAGCCGACAG CCGGGGCGGG CGACTTCCGC AGCGCCGCCC  
GTGTCAGCCT GTGGCAGGAC TTCTACCAGG CCGGCGGTGA ACCGAAGTTC  
AAGCCGGTGT CCTTCGAGCA GCCGCTGTAC ATCCTCTATT CCAGCGGCAC  
CACGGGCGTG CCCAAGTGCA TCGTCCACGG TGTGGTGGC ACCCTGCTGC  
AACACGTCAA GGAAGTGGG CTGCATACGG ACCTGACGGC CGACGACACG  
CTGTTCTACT ACACCACCTG CGGCTGGATG ATGTGGAAGT GGCTGGTCTC  
AGGGTTAGCC TTGGGCGCCA GCCTGGTGCT GTTCGACGGC TCGCCGTTCC  
AAGCAGGTGC CGAGCGCTG ATCGACCTGA TCGACGCCGA GAACATCAGC  
CTCTTCGGTA CCAGCGCAA GTTCATCGCC GCCCTGGAAA AGGCCGGCGC  
CAAGCCGCGC GAGACGCACA GGCTGCGCCG CCTGAAGGCC ATCCTCTCCA  
CCGGCTCGCC GCTGGCCAC GAGAGCTTCG AGTACGTCTA CCGCGATATC  
AAAAGCGACG TCTGCCTGTC CTCCATCTCC GGCGGCACCG ACATCGTCTC  
CTGCTTCGCC CTCGGCAACC CGACCCTGCC CGTGTGGCGC GCGAGCTGC  
AGTGCAAGGG CCTGGGCATG GATGTGCAGG TGTGGAACGA GGCCGGCCAG  
CCAGTCATCG CTGAGAAAGG CGAGCTGGTC TGCGCCCGCC ACTTCCCGTC  
GATGCCGGTC GGCTTCTGGA AGGACGCCGA TGGCGAGAAA TTCCGTAGCG  
CCTACTTCGA CACCTTCCCC GCGCTCTGGG CCCACGGCGA CTATGCCGAG  
ATCACCGAAC ACGATGGCCT GGTGATCCAC GGCCGCTCCG ACGCCGTGCT  
CAACCCCGGC GCGGTGCGCA TCGGCACTGC CGAGATCTAC CGCCAGGTGG  
AGAAGGTCGA GCAGGTGCTG GAGTCCATCG CCATCGGCCA GGACTGGGAA  
GGCGACGTGC GCGTGGTGCT GTTCGTGCGC CTGCGTGACG GCGTGGCGCT  
GAGCGACGAA CTGCAGGCAC AGATCCGCCA GGTGATCCGC GCCAACACCA  
CGCCGCGCCA TGTGCCGGCC AAAATCATCG CCGTCGCCGA CATCCCGCGC  
ACCATCAGCG GCAAGATCGT CGAGCTTGCC GTGCGCAACG TGGTGCACGG  
CAAGCCAGTG AAGAACACCG ATGCCCTGGC CAACCCGCAA GCACTTGAGC  
TGTATCGCGA TCTGCCGCAA CTGCAGTCAT GAGCCGGTAA GCGACACCGT  
AAGGGCAATG GACTGCCACT CCAGCATCTA TAGTGGATGC GCATAACCCG  
GACAGGATGT TGCTGATGCA GGTGGTTTTT TGCCTGCTGG TTGCCTTGCT  
GTATGTCGGC GCGGTGGCCG CTGACGAACC ACTGGCCTTG CATATGCCGG  
ACGCCCGGCC GCTGACCCTG TACCACGACG AGCGCGGCCA CCGCATGGTC  
GGCGACATCA CGCTCGCGGC CATCACTCTC AGCGGCCGAA CCGCACGCAT  
CGTCGACGAG CCCTGGGCCA GAGCCAGGT GAACGTCGCC AGCGGCCAGA  
ATCAACTGAT CATCCCGCTG TCGCGTACCC CGGAGCGTGA GCAACGCTAC  
ACCTGGATCG TCCCGATCAT GCCGCTGGAG CGCGCCTTCT TCAGCCTCGA  
CAAACCTGTC AGCAGCTTCG CGCAGGCACG CCAGCGCTAC CGGCGTATCT  
GCGTCGGGCT CGGCACCGCT CAAGTGGAAG TCCTGCGGCG CGAGGGTTTC  
GCCGACGAGC AGATCATCCA GCTCAAACCT GGCGAAAACC CCGCCATCCT  
GCTCGAACGC GGGCGTCTCG ATGCCTGGTT CACCGGGATT CCGGAGGCGC  
TGTACATTTG GCACAAATCT GCGGAACAGC GCCGCAAGCT TTATCAGAGC  
CCGGTCCTGG CCAGCACCGA CCTGTACCTG GCCTGCTCCA GGATCTGCGC  
CCCGCAGATC GTCGAGCAAC TGCGGGCCGC CGTGTGCAA CTGGAGGCCA  
GCGGCGTCAG CCCGCGCCTG CGCCAGGCCT ATCTACCCGA GCTCGATCGA  
CGGTGAGCAC CCAGCCACG CTGCCAGGCC TATAGAACAG ATCATCCACG  
CCGCAGGCCT GCGCCGGCTT GAGCATAACG CCAGCACTGG GGCAGACTGC  
GACTACCCCG CACAACCCCG GTGACAATAT GGACGTTGTC AAAACCCCTCA  
AGCCCGGCAA AAGCGCTTGC AAGAACGCTA CCGCGAGCAA

CTCGTCGCCG TCCGCTACCG CCTCGACCGC AAAACCAACA CCCACTACAC  
CACGGTCGAA CTCATCGTCG AACAAAAGTA CGCCCTGTAC AAAACCCCGC  
CACCCGCTCC CACACCTCCG GTAGCCCTGC GCATCTTCCG CCACGAAAAC  
GACCTCCAGC GACTGATCAG AAGCGCCGGC GGCAAGTGGG ACCGTGAGAA  
TCAGGTGTGG CTGATCGAGC GAAGCGAGGC CGAGAGGCTG GGGCTGGCGG  
AACGGATCAT CTGGACATAA TGGCTATATG TGGACATCAA GATGCCTAGT  
AATAGCCACA AACACCCAGC ATCGGACACT ATGCCTACCC CTAGGCATGT  
CGTATAAACA CTAGTTATAC AAATATCATA TGAACGACGC GACCCTAAAG  
CTAGTTAATC AAAGACAGCT CGTATCGGTA ATGAATAAAA CGAAGTGGAC  
TGAGCTGTGC AATTCATTTG ACTGCGAGAA TAAAGCATCT CCGAATGTTT  
GCTATAAATT AATTTACAGT GAACAAGAAT TCGGTTTTTC AAAAATATGG  
TGGAAATCAGC TTTTGCATGA GTGCGAAGCA ATCGAATGGA TTGATTTCAA  
ACTAGTATTG CGAGAACACC GTGGCAATCT ATTGCCAGAC AAAGAAATTG  
ATATAAGCAA ACAAATTAAG GAAGCACTAC AGGCGCATGA CATCCCTTAC  
TCTGTTGAAG GAGAAAATCT TAGGGTTTGG GGCTATATTA GCGCAGAAAA  
GAGTCCAGTA TTCGTATAAC AATTGGTTCA AGTCACTCGC TTCGCTCGCT  
CGGGACCGGC TAAAGCCGGC CCCTTAACCA AACGTTAGAT GCTTATGAAA  
AAGACAGTTC TCATACTCGT CCCAGCATTG CTACTCTCAG GATGTGGCGA  
CCCTGAATTT CACTACCAA ATGGTGACGA ATCAAAAAAT ATAACGCTAC  
GCATCCCTAA GAATTACATA AATTATTTCC CTGGCGTGAA GTACGAAAAA  
GACGGACCTG TCGTCATCAG ATTTTCATAT CCACAATTGG AGCCACTGAC  
AAAAGCCCTA CCAGAAGAGC AAAAAGTAAC TGTCAGCATT AGTCATTTAT  
CCAGCCTGGA ACTCACCACC CAAGAAACCA GAAACCCCTA CTGCGAAAAA  
GATAAAAAGT GGAAACTCCT ACAGGCGGCG GGCATTACAG GAGAGTTCTA  
TAAATTCATC GGAAAATCTC CCGGCAGCGC CAGTGCAGAT ATAACCTATA  
AGCCCATCAA AAAAACACTT GGCCTTTACT GCATTACATG CGTGGAAAAA  
GCAAATTGTG AAATTCACGC AGTATCTAGC CAAGGAATAA GCTATTCCGC  
ATTTTATACA GAAGACTTAA TGCCAGATAA GTGGCACTCT ATTTACATGG  
CAGTCGACAA AATCCTTAGC AAATTTACAG CATCGTCGAA AGGCATCTAA  
CAATTGGTTC AAGTCGCTCG CTTGCTCAC TCGGGACCGG CTAAAGCCGG  
CCCCTTAACC AAGCGTTATG CAAGCAGTCA CCCATGAGGA AAGCACCCAT  
ATGGAGCCAG TATGAAATTG AGCGACATAA GAGCTCTAAT CATTGAGTCG  
CCAGGATGGC GAACAGTATT TGCATTTATT GTCCCACTAA TCGCAGGGAT  
TCTGTCGGGA ATATTTCGTAT CAGAAATAAC GCATAGCTCC GAAATTGTTT  
GGAAGGAATT TTATAAAGCA AAAAGCTTCT ACGGGCTATT GGCTTTGAGC  
TTGTGCATGT ATTTTACAA TAAAGCCATT TATCTACATG AAAGAGAAAT  
TTCTCGCTTC CTAGACGAG ATTACTGCAC CGCTTACATG AGAAGCAAAT  
GCCTGCCAGA GGCTGCAGAG CGATACAAAA AGCTTATACG CTCTGGCGAC  
GGCGGCGAAT TGAAGCAAGC AATGGATGAA CTGAAGAAGG TGCTCAAATG  
AAAGTACTGG CCAGCCAGA TTTTAATGCA AAAGTGCCGG CACTAAACAC  
AGAAACCATT AGTAGCCTTT CTGCATTCAT ATCAAGCGCA GAGCAATATG  
AAAAAAATG CTTCATATTG AAGAATGTAA ACTCAATGTC TCTTCTTGAT  
GGCGATATAT ATAGCGCAAA AATCAACTCA AGCAGACTAT ACTTCACCAT  
CGGAGCTGAT GAGCAAGGCG ACTACTTGCT GCTATTAGAT ATAGCTGCCT  
TGCAAACCGC ACCATCTGTC AAAAGTAGTG CTTTCTTCAC AACAAACAAC  
CCAAAACCA ACAGCTCACT CAATCCGAAG CTCAACTCTG CAATCAACCC  
AAAGCTAAAC TCAGCAATAA ACCCAAATTT AAATTGCGCC ATAAATCCGA  
AGCTGAACTC AGCAATAAAC CCAAAGCTAA ACTCGGCCAT CAACCCGAAG  
CTGAATTCAG CAATAAATCC GAACTAAAC TCAGCAATAA ATCCAAAGCT  
AACTCAGCA ATAAACCCAA AACTAAACTC CTCATAAAC CCAAGGCTCA  
ATCGAAGCTA TGGCGGCCCG TATTGTACG ATGCGAACCT TAATCAAGAA  
GCGTACTCAG TTAGAGCCAA TAACAAAATC GAAATCCTGT TCAATTGCGG  
CGGAGATTTT TATGGCTTTC TTGTAAGCGC TAACGACCGA GTGAAGATTG  
AGTTCGATAC AGGAAATACC TGGACAGGTT ATTACGTTAA AGCCAATGAA  
AAAGTTTGGC TTAGATATTC GCTTAACAAC GAATGGTTAG GGCTACTTGT  
CTAGCCCGCA TAACAAGTCG CTCAAATCGC TCACTTCGTT CGCTGGGACG



GGCTAAAGCC CGCCCCCTTAG CTTATCGTTA GGCAAAAAA TAGCAGGCAG  
GCTCAGTAAT ATGAAGTTCG ATAGAATAGC TCGTGAAGCG TTTGGGCTCAG  
TGCTTGGTCC ACTGGGGTTC AGCTGTAGTG AGTCGAAGGC ATGCACCTTC  
TATAAAAAAG TCGGCACTGA GCTCTATCAT TTTGTCATGC CAGATCAATT  
AAGCGGCCAG GAAAAGTATG ATATTAAAGT TTTTTCCAC TCGCCGCTCT  
TAGAGCCCAAC CGCATGGAAT GACAAGTTTC CGGACACCTT GGGGATTCCC  
ACAGATAGCT GGAGTTATCT TTCTAGCCGT ACTGGCGTTG GTCCACGACA  
AGAGCTGTTT TGGTGTGCGA CAGAAGAAGG ATTTATGCGT AACTTTGAAT  
CAAAGGTAAA GCCCCGACTA CTTCAATTTG TAGCCCAATA TTTTGATTCT  
ATACAGACAT TGGAAGAGGC TATTCCACTA ATCAAGAGCA GGCATATGT  
GGCAGTGGCG TCTACGCTAA ATGCTAACTA AGCAATGCCA AGGTCTTCCA  
CCGGCACCTC CGTATCGGCC TTGACAGATA GCAGCAATGA GTTTCCAGCA  
AAAACCAATG CGCCGCTTGC AAGGCTGTTT CGGGTTAGCC ACAGTGCGGT  
ATTCATTACC TGCCTCCGAC TCGATACCAA TTGCCTAACA ACTGGTTCAA  
ATCGCTCGCT CCGCTCGCTG GGACCGGCGA AGCCGGCCCC TTAACCAAAC  
GTTAGGCTAC ATATGAGAAT CAGCGCAGAC CAGCTTGCTC AAGAATCACT  
GACTGAGTTC GCGGTGCTGG CGGCTAAGCT TCTGGCAACG CGAGAGCTTA  
GCCAGTTGTC CGAGAAGTTT GGGTATGCAC TGGCCTTCGG AAGGGAACCG  
GCGGCTGCCA TAGCTGAGGA CCTTGCTAGG TGCTTGTGCG GACAAAATGC  
TTCGCCGGCA TCTGAATACC CAAAATCAC CGTTAAGTAT TTCAAGGAAA  
ACGAAAGTAG TCTGTTGGCA CTCGTAGAGT GTTATGTACA AATGACCGCA  
AGCGCAAACA TTCTTTTAGA GCTGGTTGCC GCACGAAATG GAGAGGCAAT  
AAATCTGTAT CTAGAAGGCT TGAGTGTGT AGCCTAACAA TGCGCTCAAA  
CGCTCACTT CGTTCGCTGG GACCGGCGAA GCCGGCCCC TAGCTTAATC  
GTTAGGTGCC TCAGGAGGGA TCATGTCTTC CACAGAAAAC AATAGTGATG  
ACTGGCGAGA AATTCGAGCA AGAGCGGACT CTATCGCTAA TGCCATTTTC  
CTCATTTCTG GCGGGGCACT TTCACTTTCA ATCTCAGTCA TCCTCAGCAA  
CAAAAGTGCC GGGTACATCA CTGCACAAGT GGCATGTATT GCGTCCCTCG  
CTTGGTACTG CTTGCTGGCG TCACTGATTC TCTTCTTGC TCTTAAGGGG  
CATATGATTC TTCAGGCATA CCTCCTACAA TTTCGCCCAA ATTACGTCAA  
TAAACATCTT AGATTTCTTA ATGGTATAAG CTGGGCCATT GGATTAACCG  
GGTTTATTTT CTTCAATTGA GGCATGTTTC TTATGGTTCG TACCGCAATA  
CTTGCCGTG CGACCTAACA ATGCGCTCAA CTGTCGCTCA CTTGCTTCGC  
TGGACAGTCA AAAGCTGCGC TTTTGCCTGC CCGTTAGCTT AATCGTTAGC  
GGTCATAAGT ATGCAGATTA ATTTCTATAT GGCAGATGAA GATCGAAGAG  
CGTTCACGA ATACCTATAT TCTCGTGGCG CATACTCGT TCCGGAGCGT  
TGGCCAACCA GAGATATTCC CATAGTCCAG GCGGCCTCCG AAGAGGCAAG  
TGAGTGCAAA GACTTCAAGA TTTTCAAGTC TGACCTCTTC CCTCAGTCCG  
AATTTTCAGAA CAGGGCTTGG ATAACGTGGC ATGAGCCAAC GAAAAGGTTT  
TACGTTTCATG GGCCTGGAAT TCAGTATCTT GTATCGTTCA CTGATGCAAA  
TGGAATCCAT CGTGGCCGCC TCTATATGGG CCTTGTTTCT CCGCGTAGCT  
TTGTTGAGCC CCACGGGCAA TCAGTTGATT GCTACGCTGA AAACGAGAAA  
AAGTACAAGG CGTTAGAGAA TTTCCATAAG AGTTGCGCGC GCTATATACG  
CAATCACTAC CGCAAAGATG AGGGTGGTTT CTACCATGGC AAAGCAAGCG  
ATATGGCTGT TCAAACTAC GCGCTCTCAA AGACGCAGCT GTGACCGCTA  
ACAATTGCTG GCAGGCGCGA CGGCCCTGAC GGGCCGCGGC CTGAGCTCAA  
ACGTTATAAC CTACAAGGAA GACCAAAGTA TGCGCCACCT AGCAATAGCC  
CTCTTAATAA TGTTCTCTAC TCAAGTTCTC GCCGACGGCA AGAGCGAAAA  
GATAGAGAGC CTAATGAAGG CACTTGGACT AGTAGACACA TGGACACAAC  
AAATTGAACA AGGAAAAATT TACAACAGAA AGATCAGCTC TCAAATGCTG  
GATCAAATTT TATCCCAGCT GAATCCAAAT AACAGAATCT CCATGGTCTC  
CAAAAAGGCT TCAGATAATT TCATAACAAA ACTATGGCCC AGAATTCACA  
CAGAAAAAAT TGTAAGAGTT TGGGCTAGTT ACTATGGCCC AGAATTCACA  
GAGGACGAGC TTGACCAATT AATTGCATTT TATACTTCCC CTCTTGCCA  
AAAAGACATC CGTGTACTC GCAGTTCAAT GGAAAAATTC TCGAAATACT  
TCCAAGAGGC CGGGCAACCA ATACTAGAAA AGGCCACCGC AGAGTACATT

CAGGAAATGA AGCTCATCGC CAAGGAATGT AACTGTACCA AGTAGCTTAT  
AACAAATTGGT TCAAGTCGTT CGCTTCGCTC ACTGCGGGAC CGGCTAAAGC  
CGGCCCCCTTA ACCAAACGTT AGGCACTGCT ATGGCCTTGG TCGAGTACGA  
ACTGATCATC AATGCGCCCC AGACGGCTGT CTATGCCGCA TCTCAGGACT  
ATTCAGTTAG GTACCACTGG GACCCCTTCC CTGAAAAAAT TGAACCTCTA  
GGTGGTGCAA CCGAGGTAGG AATTGGGGTT AAGACACTTG TAGTCGCCAA  
GTCTGGCTTA ACAATGGAAG TCGAGTTTGT TCAGGTTGCT CCTCCTACAA  
CGGCAGCCAT AGTCATGACC AAAGGCCAG CATTTCATCAA GAGCTTTGGT  
GGTAGCTGGG TTTTCAAGCC CATCACCGCA AACTCTACAA AGGCAAAATT  
TCGCTACTCC ATAAAAACCA AGAAATGGGC AATACCCATA ATCTCAGAAT  
ACGTAGCAAG TCTTTATTTT AGAAGAGCAG TTAAGGCCAG GCTTGCCGGT  
CTTAAAAAAT ACTGCGAGCA AGGCGCCTAA CAAATGGTTC AAGTCGCTCG  
CTTCGCTCAT TCGGGACCGG CTAACGCCGG CCCCTTAGCT TAATCGTTAG  
GCTGGCCGAA GATATGAGTT ACAAGAGATG GATTTGTGTC CACTGCGATA  
CAGCCAACAC CACAGCAACA GATATTTGTT CAAAATGTCA CAGATCCAGC  
TATGAAGAGC CGGCAATAGC TGAAACTCCA ATAGCTAATT CTTACCAAGG  
CATACAGCTG TTAGGCTCTT GGCTTTTTTAT CCCACTAACC CCATCCATTA  
TGGTAATTGC AATAAGGGAT GAAGTCTGGT GGTTCGTCCC ATTTGGGATC  
GCAGTTATTG CGCTCACAAT ACTAAGTGAA AAATCTAAAT TCTTAATTTT  
CAATACTACT TGGTTCAAAA ATATAGCTTT ATTTTATACC CCAGCAGCGG  
GTGTGCTTTT CCCTCTTAGC GTTTTTCTCG GAAAGAATTG GCGGCCGCA  
TTTATGGCAA TGCACGTGGT TGTTACACCTA CATGCTGCAT TTAACATGCA  
CGCACACTGC CAAAACCACA AGCCAAGAAA TGAAAATTAA CGATAAGTAC  
TCTAGGCCTA ACAATAGGTT CAAGTCGCTC GCTTCGCTCA CTTGGGACCG  
GCTAAAGCCG GCCCCTTAAC CAAACGTTAG GTTCCATCAT GTGCTACATG  
GCCATAGTAA GTACAACATC CGAATCTGAC CTAACGGCAC TAAATACGCC  
CCTAGCGCAG TTCTCAAGGA ACGTAAACGC CATAACAGAA GCGGCCCTCC  
TGCGTTACCC AAATAAATGG TTCCTTGGCT CAAAAGACGG TTGCAGTTGT  
GCATTCAGGC ATCTTGATCA AAATGCTACA GATCTTGGTT TCTCAGAGCC  
CGTGGACTGG TGGGAAGAAG ACCAAGAAGA TATAGATGCT ACGCTTCAAG  
TAGTAGAAGC ATTCCATACG ATATTGCGCG ACGGCCATAA ACTTGACTGC  
TACATGCTT GGGCCAACGA TAGCAAGGAA CCTAAAAATC TTGCAGGTGA  
TTTGTGGTT GACCTAAATA AAGTTGGCGC CAAGAGTTTC CGTTTTTTTG  
AAGGCTATCG CTTTGAGCTC GAAGCCAGAA CCTAACTAGT GGTCAAGCC  
GCTCGCTTCG CTCACTCGGG ACCGGCTAAA GCCGGCCCCCT TAACCAAACG  
TTAGGCCCTA TATGCAGTAC TCAATTGCTG ATACTGAGAG TTTTCATCCC  
GTCATGGATG CTGAGATCAA GGGGCACTGC GAATTCCTCCG TTGACCTAAT  
TTTTGTCCCG GACATTGCAG AGTGGGCTGC TTCCAGATGC GGAAATCTCA  
TCGGCAATCC TGTGCAATG GCAGTTAGAG ACGGAGCCAC TAAGGGGGCA  
GGCATTCTCA TTAGGCAATC AATAGACGAA TCGCAAGTCG ATAGCATCCT  
ATCTCGAATG GAGTTCGGCG GCTTTGATCG TGCGCGGTCC ATACTGTCTT  
CACCCGAGAA ATTTATGCGG CATCTAGTCC TCCATGAGCT CGCGCACTTG  
ATTAACAATT GGGGGCAAGA CCGAGAAGAC GACTGTGATG AATGGGCTTT  
TAAGCGTCTT GGTGCTAGGG CCTAACAAATG CGTATATGGA CTCTCCCCAC  
AAGTAGTGGG CAAAGCTTTT GCTCTGCTCC TTTCGTGCGT GCGGTTCCAT  
GCGTATATCC GGCCTTTTCG GAGCTCAGTG CTCTGGCCAT TCTGTAGTTC  
GCGCAGCGGG GGCCAGCGC TCAAGCGATC TCCAGGATCA GCATCGTCAC  
AGGCCTGCTG TGCAGCCGGC CGAAGGAAGA CCTTCCCTCT GCTCTCCCTC  
TCCCCATCG CTCCTAATGG GCGCGTCGGC AGCCCAGGCT GCCCCAGCCT  
ACAAGGAGCA AGAACATGCC AAACCTGACG TTGCATCGCG GTGAGAAGAA  
CAACTGGTGG CCAGCACCCG AGGTGCGGCT GATCTGCGGG ATGACCCTGT  
ATGGCCCCCTG GCAGCCGACG GCGGTCAGCT CGCTCTGGAA CACACTCAAG  
AACGAGGTGA AAGGCGCCCT GGGCACCTCG CTGGAGAAGA AGGTGCGCGC  
CTACGCCCAA TACCTGCGGG CGACCGGTCC CCCCTTTGCA CTGGCGACGG  
CCTGGACCGA AGTCGGCTCC TTCACGTCCG ATTACAATA CGTGATCCGG  
ATTCCCAATG CACATTTCTT CTACTGGGGT GGCACCAAGG ATGCCCCGGA

CCTCGGCGCT	GCGGCCGCGT	GGACGACGCC	CGAGCAGGTG	ACAGCCGACT
TCATCGTGCT	CAACGCTCCC	ACCGTGGCGG	CCTCAACCAT	CCTGGGCTTC
GGGCACCACA	CCGGCACGCG	GGAGATCACC	TTCTTCCACG	ACCTGCCGAT
CGGCCTGATC	GAGTCTGTGA	ACGGCAGGCC	GATCACCAGC	TATGCGATCA
AGAGCAAGTC	CGACCTGAGC	TTCGACGAGA	GGATCAAGTA	CGCCAAATAC
CTGCGCTAGC	CGCGTACCCC	GCGTCCGAGA	GGCTTAGAAG	CTAGGGCGGC
CGGGGCTCTT	CGGGGGGGTG	TCTTCTCTGA	TTTCTCTAAG	CTTGAGTTCC
ATCGCCAGT	TGGCCGCTGC	GCCCGTGGGC	GTCGCCAACG	GTGCGGGCGT
CGGGGCGGCA	GCCTGCGGGG	CGGTGGGGTT	GTCCTTGTAC	AGCTTGCAGC
TGAGGCGCAC	GTTGTTGGCC	GAGTCGGCGT	TCTTCACTGC	CTCCTCCTCG
TCGATGACGC	CTTCATGAAC	GAGGTCGATC	AGCGCCTGGT	CGAAGGTCTG
CATGCCGAGG	TTCTTCGACT	TCTCCATGAT	CTCCTTGAGC	TCGGAGAACT
CGTTGCGCTT	GATCAGGTCG	CGTACGGTCG	GCGTGCCGAG	CATCACCTCT
ACGGCGGCGC	GGCGCTTGCC	ATCGACGGTC	TTGACCAGGC	GTTGGGAGAC
GAAGGCGCGC	AGGTTGTTGC	CGAGGTCGTT	GAGCAGCTGC	GGGCGGCGCT
CTTCGGGGAA	GAAGTTGATG	ATGCGATCCA	GCGCCTGGTT	GGCGTTGTTG
GCATGCGAGG	TGGAATGGC	CAGGTGACCG	GTGTCGGCGA	AGGCCAGGGC
GTGCTCCATG	GTTTCGCGGT	CCGGATCTCT	CGCGATCAGG	ATTACATCCG
GCGCCTGGCG	CAGAGTGTTC	TTCAGCGCGG	CGTGGAAGCT	GCGGGTGTC
ACGCCGACTT	CGCGCTGGTT	GATGATCGAC	TTCTTGTGCC	GGTGACAGTA
CTCCACCGGG	TCCTCGATGG	TGATGATGTG	GCCGCCGCTG	TTGCGGTTGC
GGTAGTCGAT	CAGCGCCGCC	AGGGAGGTCG	ACTTGCCGGA	GCCGGTACCG
CCGACGAACA	GCACCAGACC	GCGCTTCTCC	ATCACCGTCT	GCAGCAGCAC
CTCGGGCAGC	TTGAGGTCCT	CGAACTTGGG	GATGTCCATC	TTGATGTTGC
GCGCGACGAT	GGATACCTCG	TTGCGCTGCT	TGAAGATGTT	GATGCGGAAG
CGACCGCAT	TGGGCACCGA	GATGGCCAGG	TTCATCTCCA	GCTCCTTCTC
GAACTCGGCG	CGCTGCTCGG	CGTCCATCAC	GCTATTGGCG	ATGGCGGCGA
CGTCACCCGG	CTTGAGCGGC	TCCTGGCTGA	GCGGCTTGAG	CACGCCATTG
AACTTGCGC	AGGGCGGCGC	CCCGGTGGAC	AGGTAGAGT	CGGATCCGTC
CTGGCTGGAC	AGGATTTTCA	GCATCTGGGA	AAGGTCCATC	GCACGCGCTT
CCATTTGGGT	GGAGTTAACA	AGGTAGGCCA	GCTTTGCCCG	GCCGATCAGC
CTGAAAAATG	GCGCCATTCT	GATGGCGCAA	CGAATGCTGG	CACAATAGCG
CCATCGCAAA	ATGAGGACCC	CGTCATGCCC	AAAGCCATGG	CCCGCCACAT
CCTGGTGAAA	ACCGAAGCCG	AAGCCGCCGC	CCTGAAGAAA	CGTATCGCCG
CCGGCGAGGC	CTTCGATGTG	CTGGCAAAGA	AGTACTCCAC	CTGCCCCCTC
GGCAAGAAAG	GAGGCGACCT	GGGCGAGGTG	CGCCCCGGGC	AGATGGTGCG
GCGCGTGGAC	CAGGTGACTT	TCAAGAAGCC	CTTGCGGAA	GTGCACGGCC
CGGTGAAGAC	CCAGTTCGGC	TATCACCTGA	TCCAGGTGTT	CTACCCGAG
TGATCCAGCG	GCTTAGCCGG	CCCAGCCGAG	GGTAATGGCG	GCCAGCACCA
GGTAACGGCC	GGTCTTGGCC	AGGGTCACCA	GCAGCAGGAA	GCTCCACCAG
GGCTCGCGCA	TCACCCAGC	CATCAGCGTC	AGCGGGTCGC	CGATCACC GG
CGCCCAGCTC	AGCAACAGCG	ACCAGCGGCC	ATAGCGCCGA	TAGGTGTGTT
TGGCCTGCTC	CAGGCGTTGC	GCGCTCACCG	GGAACCAGCG	GCGCTCATGA
AAGCGCTCGA	TGCCACGGCC	CAGCGCCGCA	TTTCAACACC	GAGCCCAGCA
CATTGCCCGA	TACTGGCCAC	CGCAGCGCAG	ACGAACACAG	GCTGGGCGCC
ACCCAGCAAC	AGGCCAGCCA	GCAGCGCCTC	CGACTTGCA	GGGCAAGCAG
GCTGGCGGCA	CCGAAGCGAG	AAAGAAACAG	GCCGAAGTAG	ACCGAAAAGT
CGAACACAGG	TGCCATCCGG	CAAAAAGTCG	GG	

Fig. 3E  
Alignment of *Pseudomonas alcaligenes* repeat  
(PAR) elements from Contig 1

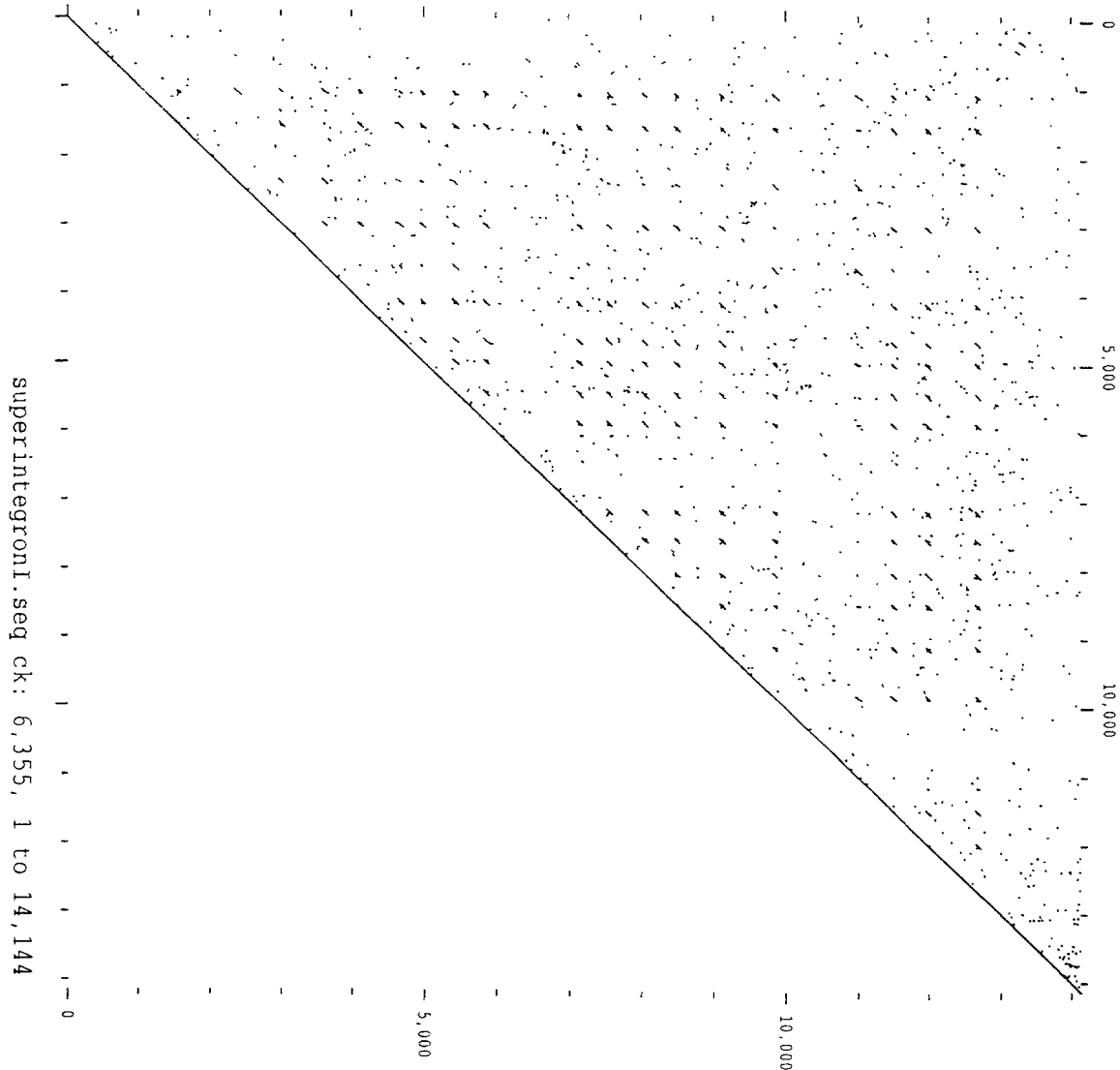
Consensus of 90% (16/18) Majority	All PARs																	
	10	20	30	40	50	60	70	80										
PAR1	GC.C	AG	A	T	G	-----	C.GCC	..G.TT	T	..A	..C	..T	..G	..C	..T	..T	..T	..A
PAR2	T.C	---	T	A	---	---	---	---	---	---	---	---	---	---	---	---	---	---
PAR3a	GC.C	CT	A	T	G-ACAG	..AAAA	..CGCTTT	T.TGCC	G	..G.TT	T	..T	..T	..T	..T	..T	..T	..T
PAR4	C	---	C.A	A	---	---	---	---	---	---	---	---	---	---	---	---	---	---
PAR5	A	---	T	A	G	---	---	---	---	---	---	---	---	---	---	---	---	---
PAR6	T.C	---	A	T	---	---	---	---	---	---	---	---	---	---	---	---	---	---
PAR7	C	---	C	A	CG AC T	..TAA	..CGCTT	..GCC	G	---	---	---	---	---	---	---	---	---
PAR8	A	---	T	A	---	---	---	---	---	---	---	---	---	---	---	---	---	---
PAR9	T	CAG	---	G	GA	..AAAC	C	G	-C GT	---	---	---	---	---	---	---	---	---
PAR10	---	C	---	---	T	---	---	---	---	---	---	---	---	---	---	---	---	---
PAR11	G--A	T	A	G	---	---	---	---	---	---	---	---	---	---	---	---	---	---
PAR12	T.C	---	T	A	---	---	---	---	---	---	---	---	---	---	---	---	---	---
PAR13	A	---	C	T	A	---	---	---	---	---	---	---	---	---	---	---	---	---
PAR14	GC.C	AG	A	T	G	---	---	---	---	---	---	---	---	---	---	---	---	---
PAR15a	GC.C	T	CT	A	T	G-ATAGT	..AAAA	..CGCTTT	T.TGTC	G	..G.TT	T	..T	..T	..T	..T	..T	..T
PAR16	C	---	A	C	---	---	---	---	---	---	---	---	---	---	---	---	---	---
PAR17	A	AC	T	A	---	---	---	---	---	---	---	---	---	---	---	---	---	---
PAR18	C	---	C	A	---	---	---	---	---	---	---	---	---	---	---	---	---	---

(SEQ ID NO:96)  
(SEQ ID NO:97)  
(SEQ ID NO:98)  
  
73 (SEQ ID NO:99)  
74 (SEQ ID NO:100)  
86 (SEQ ID NO:101)  
73 (SEQ ID NO:102)  
74 (SEQ ID NO:103)  
73 (SEQ ID NO:104)  
83 (SEQ ID NO:105)  
74 (SEQ ID NO:106)  
69 (SEQ ID NO:107)  
73 (SEQ ID NO:108)  
75 (SEQ ID NO:109)  
74 (SEQ ID NO:110)  
74 (SEQ ID NO:111)  
73 (SEQ ID NO:112)  
85 (SEQ ID NO:113)  
72 (SEQ ID NO:114)  
73 (SEQ ID NO:115)  
72 (SEQ ID NO:116)

DOTPLOT of: superintegronI.pnt Density: 16075.00 May 25, 1999 16:22

COMPARE Window: 21 Stringency: 14 Points: 52,932

superintegronI.seq ck: 6,355, 1 to 14,144



09/701 626  
Fig 4

Fig 5

Squiggle plot of: Parla.mfold May 26, 1999 11:38

(Linear) MFOLD of: Parla T: 37.0 Check: 3607 from: 1 to: 78 May 26, 1999 11:38

Length: 78 Energy: -16.4

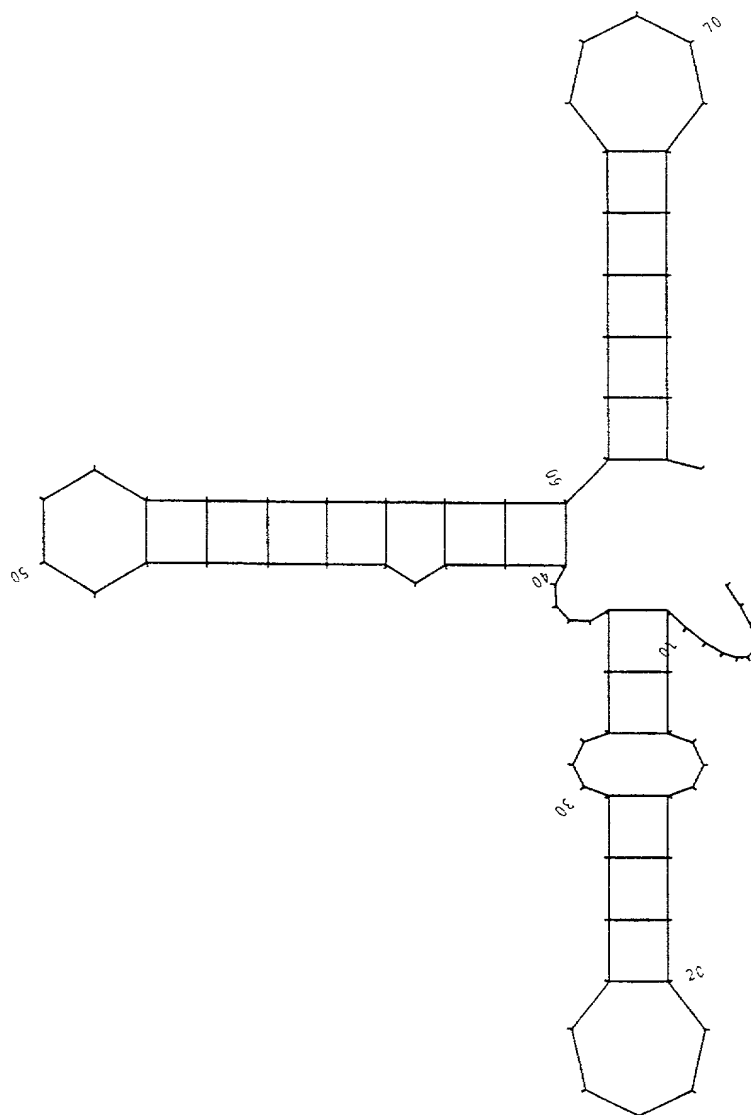


Fig. 6A

Family 1 of *Pseudomonas alcaligenes* repeat  
(PAR) elements

Identity	TAAC.....T.A.....CT.C.TC.CT.....C.G.....A.GCC.....GC.CC..T..C..CGTTA...	(SEQ ID NO:117)
90% (24/27)	CTAAC.A.TGGTTCAA.....GCTTCGCTCACT---GGGACCGGCTAAAGCCG---GCCCT-TAA-CCAAACGTTAGG	(SEQ ID NO:118)
Majority	GCCTAACAAATTGGTTCAA-GCCGTTGCTTCGCTCACT---CGGGACCGGCTAAAGCCG---GCCCT-TAA-CCAAACGTTAGGC	(SEQ ID NO:119)
PAR1	ATCTAACAAATTGGTTCAA-GTGCTCGCTTCGCTCACT---CGGGACCGGCTAAAGCCG---GCCCT-TAA-CCAAACGTTAGGT	(SEQ ID NO:5)
PAR2	ACCTAACATGGCGCTCAACCGCGCTCGCTTCGCTCGTGGACGCTGCGCATAAAGCCGCGAGCGCGGTAGCTCTACGTTAGGC	(SEQ ID NO:6)
PAR7	ACCTAACAACTGGTTCAA-GTCGTTGCTTCGCTCACT---CGGGACCGGCTAAAGCCG---GCCCT-TAA-CCAAACGTTAGGC	(SEQ ID NO:11)
PAR8	GCCTAACAAATTGGCTCAAG-TGCTTCGCTTCGCTCACT---CGGGACCGGC-GAAGCCG---GCCCT-TAG-CCAAACGTTAGGT	(SEQ ID NO:12)
PAR10	GCCTAACAAATGGTTCAA-CCGTTGCTTCGCTCACTG---GGACGGGCTAAAGCCC---GCCCT-TAA-CCAAACGTTAGGC	(SEQ ID NO:14)
PAR14	CCCTAACAAATGGTTCAAAGCCGTTGCTTCGCTCACT---CGGGACCGGCTAAAGCCG---GCCCT-TAA-CCAAACGTTAGAG	(SEQ ID NO:18)
PAR15	CTCTAACAAATGGTTCAA-GTCGCTCGCTTCGCTCACT---CGGGACCGGCTAAAGCCG---GCCCT-TAA-CCAAACGTTAGGC	(SEQ ID NO:19)
PAR18	CCCTAACTAGTGGTTCAA-GCCGCTCGCTTCGCTCACT---CGGGACCGGCTAAAGCCG---GCCCT-TAA-CCAAACGTTAGGC	(SEQ ID NO:22)
PAR19	GCCTAACAAATGGTTCAAGTGGTTC-GCTTCGCTCACTG---CGGGACCGGCTAAAGCCG---GCCCT-TAA-CCAAACGTTAGGT	(SEQ ID NO:23)
PAR20	ACCTAACAAATTGGTTCAA-GTCGTTGCTTCGCTCACT---CGGGACCGGCTAAAGCCG---GCCCT-TAA-CCAAACGTTAGGC	(SEQ ID NO:24)
PAR22	GCCTAACAAATGGTTCAA-GTCGTTGCTTCGCTCACT---CGGGACCGGCTAAAGCCG---GCCCT-TAA-CCAAACGTTAGGT	(SEQ ID NO:26)
PAR34	TTATAACAAATGGTTCAAGTGGTTC-GCTTCGCTCACTG---CGGGACCGGCTAAAGCCG---GCCCT-TAA-CCAAACGTTAGGC	(SEQ ID NO:38)
PAR36	GCCTAACAAATAGTTCAA-GTCGCTCGCTTCGCTCACT---TGGGACCGGCTAAAGCCG---GCCCT-TAA-CCAAACGTTAGGT	(SEQ ID NO:40)
PAR37	ACCTAACAAATTGGTTCAA-GCCGCTCGCTTCGCTCACT---CGGGACCGGCTAAAGCCG---GCCCT-TAA-CCAAACGTTAGGC	(SEQ ID NO:41)
PAR38	ATATACAAATGGTTCAA-CCGTTGCTGCGCTCACTG---GGACGGGCTAAAGCCC---GCCCT-TAA-CCAAACGTTAGTC	(SEQ ID NO:42)

Fig. 6B  
Family 2 of *Pseudomonas alcaligenes* repeat  
(PAR) elements

Identity	...AAC.....G.TCAA...C.CTC.C...G.TC..T.G.....G.C...GC.C...TTA.....CGTT..	(SEQ ID NO:120)
90% (13/14)	...AAC...GTTCAA...C.CTC.CT.CGCTC..T.GGGACCG..C...CG.....GC.CC...TTA.CCAACGTT..G	(SEQ ID NO:121)
Majority	GCCTAACAACTGGTTCAAG-TCGCTCGCTTCGCTCGCTCGGACCG-GC--TAAAGCCG-----GCCCG--TTAACCAAAACGTTAGGG	(SEQ ID NO:122)
PAR2	ACCTAACATGGCGCTCAACCGGCTCCCTTCGGTCGCTGGACGCTGCGGATAAAGCCGCGCAGCGCCGGTTAGCTCTACGTTAGGC	(SEQ ID NO:6)
PAR17	GCCCAACAAATGGTTCAAG-TCGCTCGCTCGCTCACTCGGACCG-GC--TAAAGCCG-----GCCCG--TTAACCAAAACGTTAGGG	(SEQ ID NO:21)
PAR21	GCCTAACAACTGGTTCAAG-TCGCTCGCTTCGCTCGCTCGGACCG-GC--ATAGC-CG-----GCCCG--TTAACCAAAACGTTAGGT	(SEQ ID NO:25)
PAR28	GCCTAACAACTGGTTCAAG-CCACTCGCTTCGCTCGCTCGGACCG-C--GTACCCG--GCCCG--TTAACCAAAACGTTGGGC	(SEQ ID NO:32)
PAR29	GCCCAACAAACGGTTCAAGACCGCTCGCTTCGCTCGCTCGGACCG-GC--TAAAGCCG-----GCCCG--TTAACCAAAACGTTAGGG	(SEQ ID NO:33)
PAR30	GCCTAACAACTGGTTCAA-ATCGCTCGCTCGCTCGCTCGGACCG-GC--GAAGCCG-----GCCCG--TTAACCAAAACGTTAGGC	(SEQ ID NO:34)
PAR44	CCCTAACAAATGGTTCAAG-TCGCTCGCTTCGCTCGCTCGGACCG-GC--TAAAGCCG-----GCCCG--TTAACCAAAACGTTAGAG	(SEQ ID NO:48)
PAR53	GCCTAACAACTGGTTCAAG-CCACTCACTTCGCTCGCTCGGACCG-C--GTTCCGCG-----GCCCG--TTAACCAAAACGTTGGGC	(SEQ ID NO:57)
PAR56	GCCTAACTACTGGTTCAAG-TCGCTCGCTTCGCTCGCTCGGACCG-C--GTTCCGCG-----GCCCG--TTAACCAAAACGTTAGGC	(SEQ ID NO:60)
PAR61	ATCTAACATGTGGTTCAAG-CCGCTCGCTTCGCTCACTCGGACCG-GC--TAAAGCCG-----GCCCG--TTAACCAAAACGTTAGAG	(SEQ ID NO:65)
PAR62	CTCTAACAAATGGTTCAAG-CCGCTCGCTTCGCTCGCTCGGATCG-GC--GAAGCCG-----GCACC--TTAACCAAAACGTTAGAG	(SEQ ID NO:66)
PAR67	GCCTAACAACTGGTTCAA-ATCGCTCGCTTCGCTCGCTCGGACCG-GC--ATAGCCG-----GCCCG--TTAACCAAGCGTTAGAT	(SEQ ID NO:71)
PAR69	GTTTAACTACTGGTTCAAG-CCGCTCGCTTCGCTCACTCGGACCG-GC--TAAATTCG-----GCCCG--TTAGGCAAAACGTTAACT	(SEQ ID NO:73)
PAR71	GTATAACAAATGGTTCAAG-TCGCTCGCTTCGCTCGCTCGGACCG-GC--TAAAGCCG-----GCCCG--TTAACCAAAACGTTAGAT	(SEQ ID NO:75)



Fig. 6C  
Family 3 of *Pseudomonas alcaligenes* repeat  
(PAR) elements

Identity 90% (13/15) Majority	..AACA...CGCTAAC...CGCTCCT..CGTCGCTGGA...GC....GCGCTAGCT..ACGTTA... ..CTAACAATGCGCTCACT..CGCTCACTTGGTTCGCTGGAC...CAAAAGCT..CGCTTTTG...GCCGCTTAGCTTAATCGTTA... ACCTAACAAATGCGCTCACTGTCGCTCACTTGGTTCGCTGGACGCAAAAGCTGCGCTTTTGCTGCTGCGCTTAGCTTAATCGTTAGGC	(SEQ ID NO:1) (SEQ ID NO:1) (SEQ ID NO:125)
PAR5	ACCTAACAAATGCGCTCACTGCGCTCACTTGGTTCGCTGGACAGTCAAAAGCTGCGCTTTTGCTGCTGCGCTTAGCTTAATCGTTAGGC	(SEQ ID NO:9)
PAR9	ACCTAACAAATGCGCTCACTGCGCTCACTTGGTTCGCTGGACAGTCAAAAGCTGCGCTTTTGCTGCTGCGCTTAGCTTAATCGTTAGGC	(SEQ ID NO:13)
PAR32	ACCTAACAAATGCGCTCACTGTCGCTCACTTGGTTCGCTGGACAGTCAAAAGCTGCGCTTTTGCTGCTGCGCTTAGCTTAATCGTTAGCG	(SEQ ID NO:36)
PAR41	GCCTAACAAATGCGCTCACTGCGCTCACTTGGTTCGCTGGACAGTCAAAAGCTGCGCTTTTGCTGCTGCGCTTAGCTTAATCGTTAGAG	(SEQ ID NO:45)
PAR25	ACATAACAATGCGCTCACTGCGCTCACTTGGTTCGCTGGACAGTCAAAAGCTGCGCTTTTGCTGCTGCGCTTAGCTTAATCGTTAGGC	(SEQ ID NO:29)
PAR13	ATTTAACAATGCGCTCACTGCGCTCACTTGGTTCGCTGGACAGTCAAAAGCTGCGCTTTTGCTGCTGCGCTTAGCTTAATCGTTAGGG	(SEQ ID NO:17)
PAR52	GGCTAACAAATGCGCTCACTGCGCTCACTTGGTTCGCTGGACAGTCAAAAGCTGCGCTTTTGCTGCTGCGCTTAGCTTAATCGTTAGGA	(SEQ ID NO:56)
PAR66	TGCTAACAAATGCGCTCACTGCGCTCACTTGGTTCGCTGGACAGTCAAAAGCTGCGCTTTTGCTGCTGCGCTTAGCTTAATCGTTAGGC	(SEQ ID NO:70)
PAR26	GCCTAACAAATGCGCTCACTGCGCTCACTTGGTTCGCTGGACAGTCAAAAGCTGCGCTTTTGCTGCTGCGCTTAGCTTAATCGTTAGGC	(SEQ ID NO:30)
PAR24	GCCTAACAAATGCGCTCACTGCGCTCACTTGGTTCGCTGGACAGTCAAAAGCTGCGCTTTTGCTGCTGCGCTTAGCTTAATCGTTAGGT	(SEQ ID NO:28)
PAR27	GCCTAACAAATGCGCTCACTGCGCTCACTTGGTTCGCTGGACAGTCAAAAGCTGCGCTTTTGCTGCTGCGCTTAGCTTAATCGTTAGGC	(SEQ ID NO:31)
PAR54	GCCTAACAAATGCGCTCACTGCGCTCACTTGGTTCGCTGGACAGTCAAAAGCTGCGCTTTTGCTGCTGCGCTTAGCTTAATCGTTAGGC	(SEQ ID NO:58)
PAR4	ACCTAACAAATGCGCTCACTGCGCTCACTTGGTTCGCTGGACAGTCAAAAGCTGCGCTTTTGCTGCTGCGCTTAGCTTAATCGTTAGGT	(SEQ ID NO:8)
PAR11	CTCTAACAAATGCGCTCACTGCGCTCACTTGGTTCGCTGGACAGTCAAAAGCTGCGCTTTTGCTGCTGCGCTTAGCTTAATCGTTAGGT	(SEQ ID NO:15)
PAR2	ACCTAACAAATGCGCTCAACCG-CGCTCCCTTCGCTGCGCTGCGCTGCGATAAAGCCGCGCAGCGCGCTTAGCTTA-CGTTAGGC	(SEQ ID NO:6)

Fig. 6D  
Family 4 of *Pseudomonas alcaligenes* repeat  
(PAR) elements

Identity	TAACA...CGCTCAA...CG.TC.CT.CG.T.C.CTGG...G.GC...AAGCC...GC.CC...TTAGCT...CGTTA.	(SEQ ID NO:126)
90% (7/8)	TAACAA.GCGCTCAA...GCTC.CTTGG.T-CGCTGGG-AC.G-GC--..AAGCC...GCCCC--TTAGCTTAATCGTTAG...	(SEQ ID NO:127)
Majority	GCCTAACAAATGCGCTCAA...GCGCTCACTTCGTT-CGCTGGG-ACCG-GC--TAAAGCCG...GCCCC--TTAGCTTAATCGTTAGGT	(SEQ ID NO:128)
	10 20 30 40 50 60 70 80	
PAR6	GCCTAACAAATGCGCTCAA...GCGCTCACTTCGTT-CGCTGGG-ACCG-GC--GAAGCCG...GCCCC--TTAGCTTAATCGTTAGGT	(SEQ ID NO:10)
PAR31	GCCTAACAAATGCGCTCAA...GCGCTCACTTCGTT-CGCTGGG-ACCG-GC--GAAGCCG...GCCCC--TTAGCTTAATCGTTAGGT	(SEQ ID NO:35)
PAR50	GCCTAACAAATGCGCTCAA...GCGCTCACTTCGTT-CGCTGGG-ACCG-GC--GAAGCCG...GCCCC--TTAGCTTAATCGTTAGAA	(SEQ ID NO:54)
PAR65	GTCTAACAAATGCGCTCAA...GCGCTCACTTCGTT-CGCTGGG-ATCG-GC--TAAAGCCG...GCCCC--TTAGCTTAATCGTTAGCA	(SEQ ID NO:69)
PAR42	CTCTAACAAATGCGCTCAA...TCGCTCACTACGTT-CGCTGGG-ACCG-GC--TAAAGCCG...GCCCC--TTAGCTTAATCGTTAGAG	(SEQ ID NO:46)
PAR72	GCATAACAAGTCGCTCAA...TCGCTCACTTCGTTTCGCTGGG-ACCG-GC--TAAAGCCG...GCCCC--TTAGCTTA-TGCTTAGGC	(SEQ ID NO:76)
PAR12	TTATAACAATGCGCTCAA...TCGTTTCGCTTCGCT-CAGTGGG-ACCG-GC--TAAAGCCG...GCCCC--TTAGCTTAATCGTTAAAT	(SEQ ID NO:16)
PAR2	ACCTAACATGCGCTCAACCGGCTCCCTTCGGT-CGCTGGACGCTGCGCGATTAAGCCGCGCGGTTAGCTCTA-CGTTAGGC	(SEQ ID NO:6)

Fig 7A. PAR-specific oligonucleotide (bottom) aligned with PAR majority consensus (top)

majority	GCCTAACAAATTGGTTCAAG-GTCGCTCGCTTCGCTCACT-CGGGACCGGCTAAAGC-----CGGCC--CC-TTAA-CCAAACGTTAGGC (SEQ ID NO:129)
oligo 1	5' TCGCTTCGCTCACTGCGGACCGGCTAAAGC-----CGGCC--CC-TTAA-CCAAACGTTA (SEQ ID NO:79)
2	5' TAACAATTGGTTCAAG--TCGTTTCGCTTCGCTCACTGCGGACCG (SEQ ID NO:80)
3	5' TAACATATTCAGTCAAGCGGA--CGCAAAACCCGCTGCGGGTCTT (SEQ ID NO:81)
4	5' TAACAATGCGCTCAACTGCG-CTCACTTCGTTTCGCTGGACAGCC (SEQ ID NO:82)
5	5' TAACAAGTCGCTCAACTGCGGCTCACT-CGTTTCGCTGGACAGCC (SEQ ID NO:83)

Fig. 7B. PCR primers for PAR fingerprints

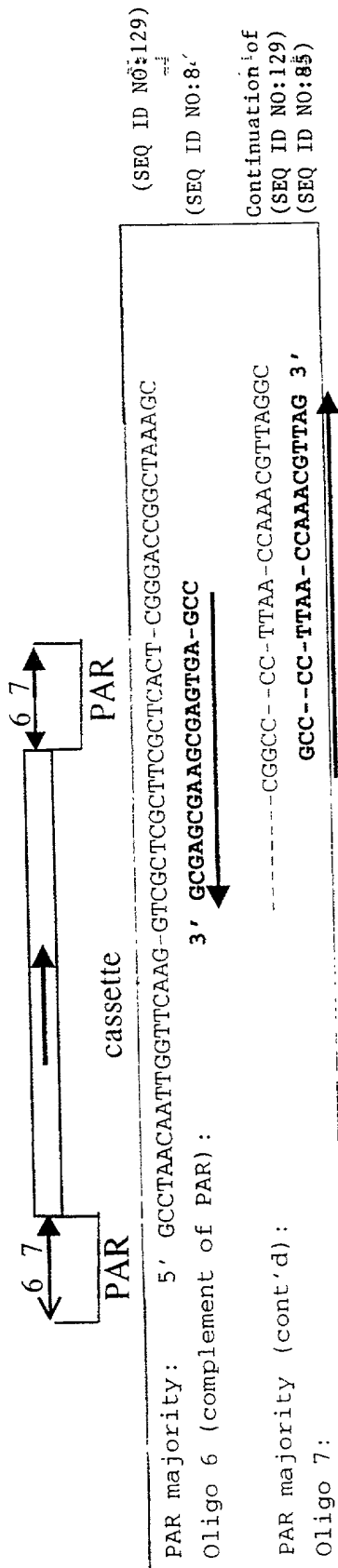
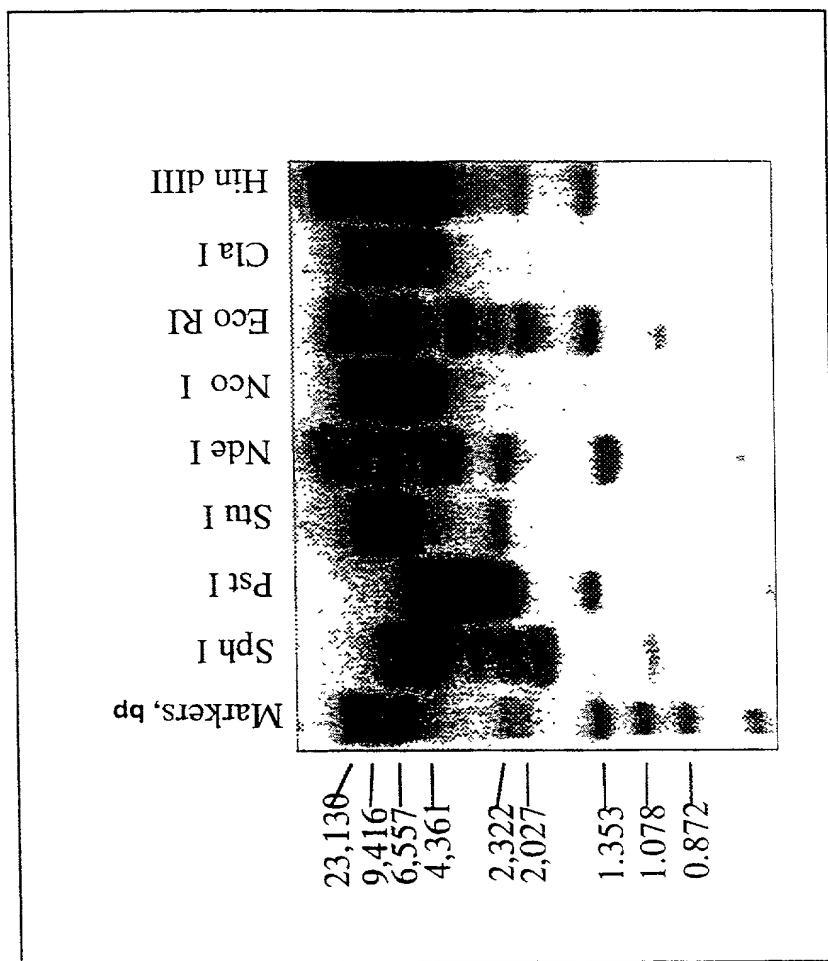
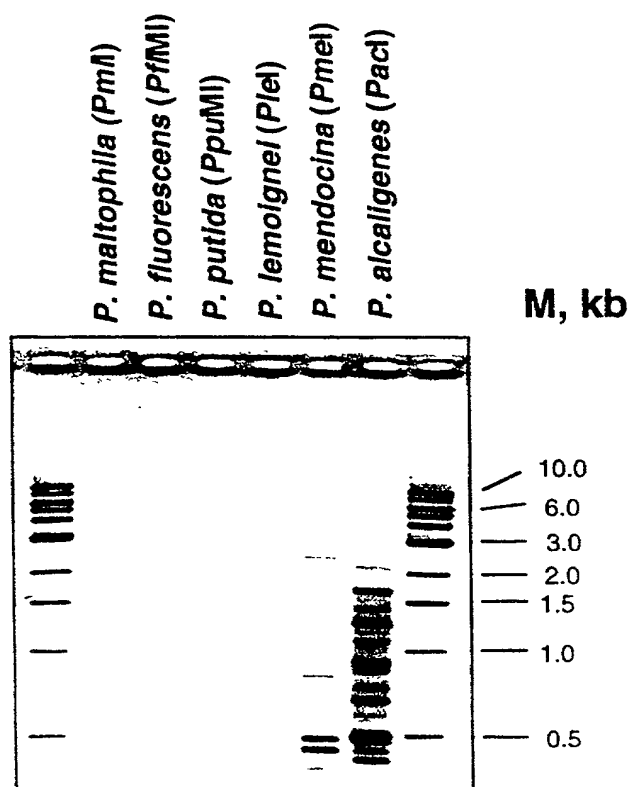


Fig 8. Hybridization of a PAR-specific Oligonucleotide 1 to *Pseudomonas alcaligenes* chromosomal DNA.

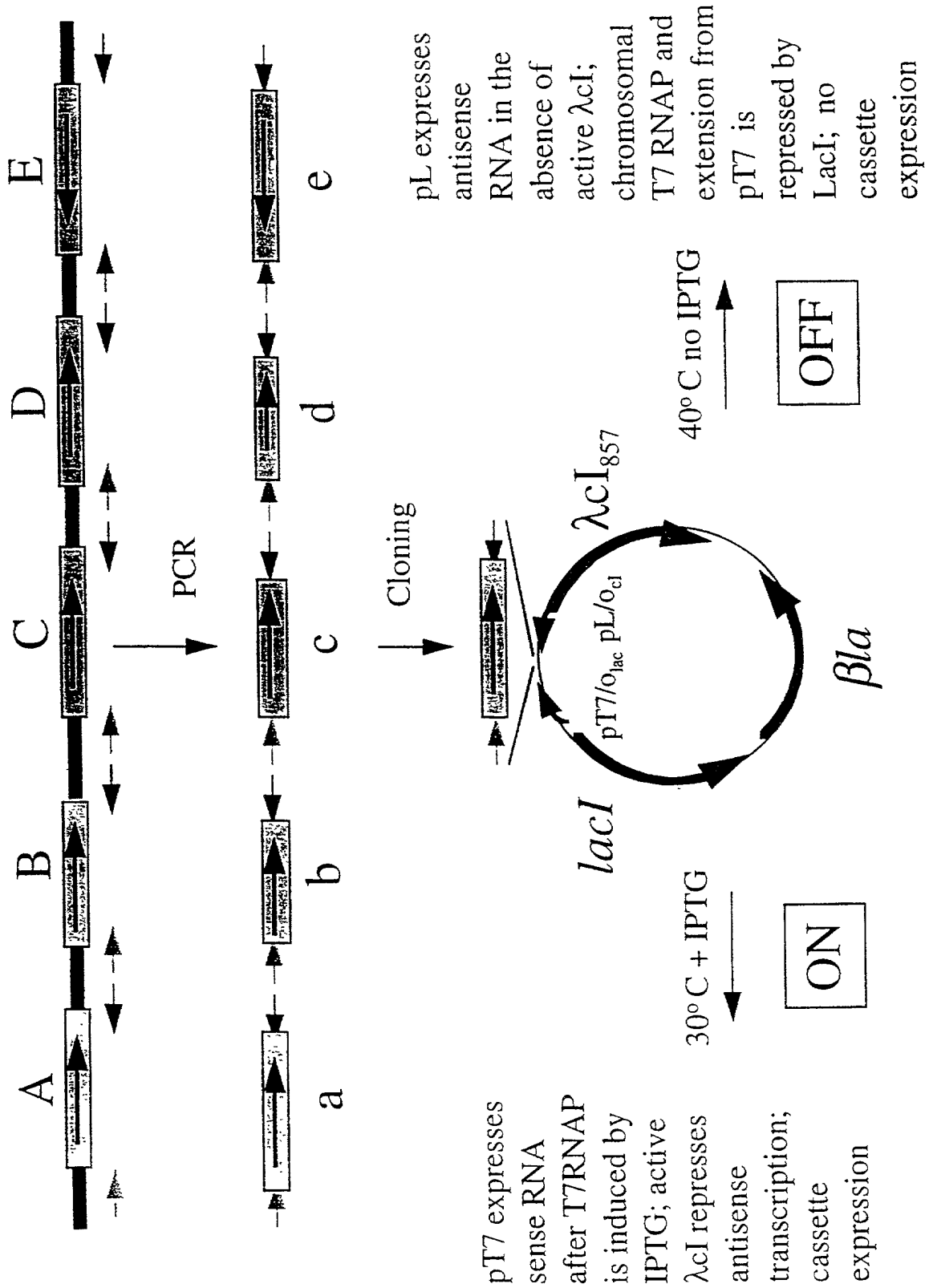


**Fig.9. Distribution of PAR Cassettes Among  
*Pseudomonas* Species**



10

Fig. . Amplification and cloning strategy



New England Biolabs, Inc.

32 Tozer Road

Beverly, MA 01915

**DECLARATION  
AND POWER OF ATTORNEY  
Original Application**

NEB-165-PUS  
Attorney Docket No. \_\_\_\_\_

As a below named inventor, I hereby declare that:

My residence, post address and citizenship are as stated below next to my name

I believe that I am the original, first and sole inventor (in only one name is listed at 201 below) or an original, first and joint inventor (if plural names are listed at 201-203 below) of the subject matter which is claimed and which a patent is sought on the invention entitled:

RESTRICTION ENZYME GENE DISCOVERY METHOD

which is described and claimed in:

☒ the attached specification or ☐ the specification in Application Serial No. \_\_\_\_\_ filed \_\_\_\_\_  
(for declaration not accompanying application)

And was amended on \_\_\_\_\_  
if applicable

I hereby state that I have reviewed and understand the contents of the above identified specification, including the claims, as amended by any amendments referred to above. I acknowledge the duty to disclose information which is material to the examination of this application in accordance with Title 37, Code of Federal Regulations, §1.56(a).

I hereby claim foreign priority benefits under Title 35, United States Code, §119 of any foreign application(s) for patent or inventor's certificate listed below and have also identified below any foreign application for patent or inventor's certificate listed below and have also identified below any foreign application for patent or inventor's certificate having a filing date before that of the application on which priority is claimed:

FOREIGN APPLICATION(S) IF ANY, FILED WITHIN 12 MONTHS PRIOR TO THE FILING DATE OF THIS APPLICATION

COUNTRY	APPLICATION	DATE OF FILING (day, month, year)	PRIORITY CLAIMED UNDER 35 U.S.C. 119
			YES NO
			YES NO

ALL FOREIGN APPLICATION(S) IF ANY, FILED MORE THAN 12 MONTHS PRIOR TO THE FILING DATE OF THIS APPLICATION

COUNTRY	APPLICATION	(day, month, year)	PRIORITY CLAIMED UNDER 35 U.S.C. 119

I hereby claim the benefit under Title 35, United States Code §120 of any United States application(s) listed below and, insofar as the subject matter of each of the claims of this application is not disclosed in the prior United States application in the manner provided by the first paragraph of Title 35, United States Code §112, I acknowledge the duty to disclose material information as defined in Title 37, Code of Federal Regulations, §1.56(a) which occurred between the filing date of the prior application and the national or PCT international filing date of this application:

Application Serial No.	Filing Date	Status (Patented, Pending, Abandoned)
PCT/US99/13295	11 June 1999	Pending

[illegible]

Gregory D. Williams  
(Registration No. 30901)

**DIRECT TELEPHONE CALLS TO:**

Gregory D. Williams  
General Counsel  
New England Biolabs, Inc.  
Tele: (978) 927-5054; Ext. 292  
Fax: (978) 927-1705

2	Full Name of Inventor	Last Name <u>Raleigh</u>	First Name Elisabeth	Middle Name A.
0	Residence & Citizenship	City <u>Somerville</u> <i>MA</i>	State/Foreign Country Massachusetts	Citizenship US
1	Post Office Address	Post Office Address 32 Barton Street	City/State/Country Somerville, MA	Zip Code 02144
2	Full Name of Inventor	Last Name <i>200</i> <u>Vaisvila</u>	First Name <u>Romualdas</u>	Middle Name
0	Residence & Citizenship	City <u>Rockport</u> <i>MA</i>	State/Foreign Country Massachusetts	Citizenship Lithuania
2	Post Office Address	Post Office Address 507 Sandy Bay Terrace	City/State/Country Rockport, MA	Zip Code 01966
2	Full Name of Inventor	Last Name <i>300</i> <u>Morgan</u>	First Name <u>Richard</u>	Middle Name D.
0	Residence & Citizenship	City <u>Middleton</u> <i>MA</i>	State/Foreign Country Massachusetts	Citizenship US
3	Post Office Address	Post Office Address 31 Donovan's Way	City/State/Country Middleton, MA	Zip Code 01949
2	Full Name of Inventor	Last Name	First Name	Middle Name
0	Residence & Citizenship	City	State/Foreign Country	Citizenship
4	Post Office Address	Post Office Address	City/State/Country	Zip Code
2	Full Name of Inventor	Last Name	First Name	Middle Name
0	Residence & Citizenship	City	State/Foreign Country	Citizenship
5	Post Office Address	Post Office Address	City/State/Country	Zip Code



[illegible]

I hereby further declare that all statements made herein of my own knowledge are true and that all statements made on information and belief are believed to be true and further that these statements were made with the knowledge that willful statements and the like so made are punishable by fine or imprisonment, or both, under Section 1001 of Title 18 of the United States Code, and that such willful false statements may jeopardize the validity of the application or any patent issued thereon.

Signature of Inventor 201	Date
<i>[Signature]</i>	30 November 2000
Signature of Inventor 202	Date
<i>P. Nair</i>	30 November 2000
Signature of Inventor 203	Date
<i>[Signature]</i>	30 November 2000
Signature of Inventor 204	Date
Signature of Inventor 205	Date
Signature of Inventor 206	Date
Signature of Inventor 207	Date
Signature of Inventor 208	Date
Signature of Inventor 209	Date